

# MEASURE WORKSHEET

This document summarizes the evaluation of the measure as it progresses through NQF's Consensus Development Process (CDP). The information submitted by measure developers/stewards is included after the Brief Measure Information, Preliminary Analysis, and Pre-meeting Public and Member Comments sections.

#### To navigate the links in the worksheet: Click to go to the link. ALT + LEFT ARROW to return

Purple text represents the responses from measure developers.

**Red** text denotes developer information that has changed since the last measure evaluation review.

# **Brief Measure Information**

#### NQF #: 0177

**Corresponding Measures:** 

De.2. Measure Title: Improvement in pain interfering with activity

Co.1.1. Measure Steward: Centers for Medicare & Medicaid Services

**De.3. Brief Description of Measure:** The percentage of home health episodes of care during which the frequency of the patient's pain when moving around improved.

**1b.1. Developer Rationale:** Both acute and chronic pain have been identified as areas requiring health care provider intervention. Many patients who receive home health care experience pain, which can interfere with activity and affect virtually all aspects of a patient's daily life, as well as impact other aspects of health status. Appropriate evaluation and management of pain is recognized as very important to the wellbeing of patients. Clinical practice guidelines identify effective interventions for chronic pain including both pharmacologic and nonpharmacologic.

High-quality home health care appropriately evaluates and manages pain, and reduction in pain can be considered a marker of high-value care for patients with pain.

**S.4. Numerator Statement:** The number of home health episodes of care where the value recorded on the discharge assessment indicates less frequent pain at discharge than at start (or resumption) of care.

**S.6. Denominator Statement:** Number of home heath episodes of care ending with a discharge during the reporting period, other than those covered by generic or measure- specific exclusions.

**S.8. Denominator Exclusions:** All home health episodes where there is no pain reported at the start (or resumption) of care assessment, or the patient is non-responsive, or the episode of care ended in transfer to inpatient facility or death at home, or the episodes is covered by one of the generic exclusions.

#### De.1. Measure Type: Outcome

S.17. Data Source: Electronic Health Data

S.20. Level of Analysis: Facility

IF Endorsement Maintenance – Original Endorsement Date: Mar 31, 2009 Most Recent Endorsement Date: Jul 07, 2015

IF this measure is included in a composite, NQF Composite#/title:

IF this measure is paired/grouped, NQF#/title:

# **Preliminary Analysis: Maintenance of Endorsement**

To maintain NQF endorsement endorsed measures are evaluated periodically to ensure that the measures still meets the NQF endorsement criteria ("maintenance"). The emphasis for maintaining endorsement is focused on how effective the measure is for promoting improvements in quality. Endorsed measures should have some experience from the field to inform the evaluation. The emphasis for maintaining endorsement is noted for each criterion.

# Criteria 1: Importance to Measure and Report

#### 1a. Evidence

# Maintenance measures – less emphasis on evidence unless there is new information or change in evidence since the prior evaluation.

**<u>1a. Evidence.</u>** The evidence requirements for a health outcome measure include providing empirical data that demonstrate a relationship between the outcome and at least one healthcare structure, process, intervention, or service; if these data not available, data demonstrating wide variation in performance, assuming the data are from a robust number of providers and results are not subject to systematic bias. For measures derived from patient report, evidence also should demonstrate that the target population values the measured outcome, process, or structure and finds it meaningful.

#### Summary of prior evaluation in 2015

The developer provided four studies that support the need to measure pain in an adult home health care population. The <u>literature</u> also suggests that home care interventions may be helpful in helping patients to manage pain. Specifically, Bach, et al. (2013) found in a small study that a particular physical therapy cognitive-behavioral intervention for pain management was effective in relieving pain.

NOTE: the 2015 NQF evidence criteria for outcome measures was to provide a rationale supports the relationship of the health outcome to processes or structures of care.

The 2015 committee noted that pain management is a significant health issue related to functional outcomes and there is definitely a relationship between the measured outcome and healthcare action supported by the rationale. Similar to all measures addressing improvements in ADLs, the Committee had a major concern about the requirement for CMS to not require improvement in function as a condition of coverage in home health, and applied the same remarks from the discussion on 0167 to all ADL improvement measures.

#### Changes to evidence from last evaluation

# $\Box$ The developer attests that there have been no changes in the evidence since the measure was last evaluated.

#### ☑ The developer provided updated evidence for this measure:

The developer provided <u>additional literature</u> that suggest gaps in pain assessment and management, as well as racial disparities in experience of pain. The evidence suggests that pain is a serious problem with adverse impact on a wide range of outcomes from functional capacity, to quality of life and mortality. The cited literature also includes examples of non-pharmacological interventions (e.g., chair yoga) that may have a positive effect on pain management in older adults.

#### Question for the Committee:

• The evidence provided by the developer is updated, directionally the same, and somewhat stronger compared to that for the previous NQF review. Does the Committee agree there is no need for repeat discussion and vote on Evidence?

#### **Guidance from the Evidence Algorithm**

Measure assesses a health outcome (Box 1)  $\rightarrow$  The relationship between the outcome and the intervention demonstrated by performance data (Box 2)  $\rightarrow$  Pass

Preliminary rating for evidence: 🛛 Pass 🗆 No Pass

#### 1b. Gap in Care/Opportunity for Improvement and 1b. Disparities

#### Maintenance measures - increased emphasis on gap and variation

**<u>1b. Performance Gap.</u>** The performance gap requirements include demonstrating quality problems and opportunity for improvement.

The developer provided performance data from June 2010 through 2016, which indicate opprotunity for improvement.

Risk Adjusted Home Health Agency (HHA) Level Performance on Improvement in Pain Interfering with Activity by Calendar Year:

		Average								
Calendar	Number	Episodes	HHA	Std.		25th	50th	75th		
Year	of HHAs	per HHA	Average	Dev.	Min	%ile	%ile	%ile	Max	IQR*
HHAs with	>=1 Valid	Episode								
2010	10,831	236	61.0%	21.0%	0.0%	51.4%	63.3%	73.2%	100.0%	21.8%
2011	11,487	278	61.9%	20.6%	0.0%	51.8%	63.8%	74.4%	100.0%	22.6%
2012	11,748	273	61.9%	21.4%	0.0%	51.4%	64.3%	75.0%	100.0%	23.6%
2013	11,893	282	62.2%	21.9%	0.0%	51.0%	64.4%	76.1%	100.0%	25.1%
2014	11,832	300	61.6%	22.3%	0.0%	50.3%	64.2%	75.9%	100.0%	25.6%
2015	11,527	335	62.6%	22.9%	0.0%	50.7%	66.5%	77.5%	100.0%	26.8%
2016	11,166	347	65.6%	23.7%	0.0%	54.5%	70.0%	82.1%	100.0%	27.6%

#### Disparities

The developer provided data tables showing disparities in performance by race, age, gender, agency size, region, disability status and dual eligible status.

Observed and Predicted Episode-Level Measure Performance by Population Group:

Demulation Crown		2016	2016
Population Group		Observed	Predicted
All Episodes		74.4%	67.3%
Condor	Male	65.5%	67.1%
Gender	Female	65.5%	67.4%
	White	65.9%	67.8%
Data	Black	63.6%	65.2%
Race	Hispanic	64.6%	66.1%
	Male Male Male Male Male Male Male Male	65.4%	67.1%
	Under 65	61.3%	62.9%
A = 0	65-74	67.8%	69.6%
Age	75-84	66.3%	68.1%
	85 and Over	65.0%	66.8%

Population Group		2016	2016
		Observed	Predicted
Disability Status	No	66.1%	67.9%
Disability Status	Yes	63.5%	65.3%
Dual Enrollment in	No	66.0%	67.8%
Medicare and Medicaid	Yes	64.0%	65.7%
	Small	62.2%	63.3%
Agency Size	Medium	64.4%	66.3%
	Large	65.8%	67.6%
	Northeast	65.4%	67.1%
Canava Dagion	Midwest	66.0%	67.7%
Census Region	South	65.2%	67.1%
	West	65.6%	67.4%

#### Questions for the Committee:

- Does the measure demonstrate a quality problem related to home health care interventions and improvement in pain interfering with activity?
- Is a national performance measure still warranted?
- Are you aware of evidence that other disparities exist in this area of healthcare?

Preliminary rating for opportunity for improvement: 🛛 High 🗌 Moderate 🗌 Low 🗋 Insufficient

#### **Committee Pre-evaluation Comments:**

#### Criteria 1: Importance to Measure and Report (including 1a, 1b, 1c)

#### 1a. Evidence

Comments:

- Evidence shows direct relation between measured outcome and intervention.
- I appreciate the developer submitting additional literature as evidence.
- There are data presented showing how pain adversely impacts functional improvement and put people at a higher risk for fales, cognitive decline and depressive symptoms. There are also data in the literature citing interventions to improve pain. I'm not aware of new studies that change the evidence base.
- Sufficient evidence is presented- both tangential and directly related. Not aware of any new evidence.
- The developer has provided additional literature that suggests gaps in pain assessment and pain management. The evidence shows that pain is a serious problem causing adverse impact on functioning and quality of life. Non-pharmacologic interventions may have a positive effect on pain management.
- Evidence relates to the outcome being measured and provides evidence regarding the reality of pain being an issue for people living in their homes and that it can reduce mobility and independence. I am not aware of any new evidence.

#### 1b. Performance Gap

#### Comments:

- Yes, there is still a high performance gap.
- The statistics are strongly suggestive of important gaps in care and lots of opportunity for improving the quality of care.
- Data were provided that demonstrate significant room for improvement in improving pain dating back to 2010. They present data from the from 2016 based on gender, race, age, disbility status, dual enrollment in medicare and medicaid, agency size, and census region. There may be some differences withing groups

(e.g. age) but no statisical test were done. The subgroup differences based on age are as one might predict.

- Performance data were provided which demonstrates a sufficient gap to warrant a national performance measure. Disparities were demonstrated.
- The data presented represents 2010 through 2016 and indicates a need for improvement. In the literature it was suggested that they are racial disparities in the experience of pain. There are also disparities in performance based on race, age, gender, disability status, agency size and region.
- Current performance data was provided and shows there is room for improvement. Slight disparities indicated when looking at the population subgroups, espcially for people with a disability, those dually enrolled in Medicare and Medicaid, and agencies small in size.

# Criteria 2: Scientific Acceptability of Measure Properties

2a. Reliability: Specifications and Testing

2b. Validity: Testing; Exclusions; Risk-Adjustment; Meaningful Differences; Comparability Missing Data

#### 2c. For composite measures: empirical analysis support composite approach

#### Reliability

**<u>2a1. Specifications</u>** requires the measure, as specified, to produce consistent (reliable) and credible (valid) results about the quality of care when implemented. For maintenance measures – no change in emphasis – specifications should be evaluated the same as with new measures.

<u>2a2. Reliability testing</u> demonstrates if the measure data elements are repeatable, producing the same results a high proportion of the time when assessed in the same population in the same time period and/or that the measure score is precise enough to distinguish differences in performance across providers. For maintenance measures – less emphasis if no new testing data provided.

#### Validity

**<u>2b2. Validity testing</u>** should demonstrate the measure data elements are correct and/or the measure score correctly reflects the quality of care provided, adequately identifying differences in quality. For maintenance measures – less emphasis if no new testing data provided.

**2b2-2b6.** Potential threats to validity should be assessed/addressed.

#### Composite measures only:

**<u>2d. Empirical analysis to support composite construction</u>. Empirical analysis should demonstrate that the component measures add value to the composite and that the aggregation and weighting rules are consistent with the quality construct.** 

#### Complex measure evaluated by Scientific Methods Panel? $\boxtimes$ Yes $\square$ No

Evaluators: David Nerenz, Sam Simon, John Bott, Zhenqiu Lin, Joe Kunisch

Methods Panelists' Combined Preliminary Analysis

#### **Evaluation of Reliability and Validity:**

<u>Reliability</u>

- Reliability testing was conducted at the both the data element and measure score levels.
- Testing of the data elements
  - Developers conducted an inter-rater reliability (IRR) analysis among nurses and physical therapists using a linear weighted kappa statistic. Testing of OASIS-C2 item M1242 was done using 2016-2017 data from home health patients in 4 states.

- Start Of Care/Resumption Of Care: kappa=0.45 (n=105 patients) ["Moderate" agreement, according to the Landis and Koch classification system]
- Discharge: kappa=0.53 (n=84 patients) ["Moderate" agreement, according to the Landis and Koch classification system]
- Testing of the measure score
  - Developers used two approaches to assess reliability of the measure score: a signal-to-noise analysis using the Adams beta-binomial method and a split-sample analysis using ICC(2,1) and ICC(3,1) statistics. CY2016 data were used in testing.
    - Signal-to-noise reliability estimates: Mean=0.95; minimum=0.74; 10<sup>th</sup> percentile = 0.87; median =0.97; 90<sup>th</sup> percentile =1.00
    - Split sample reliability estimates: IRR(2,1)= 0.90; IRR(3,1)= 0.90 [NOTE that testing data limited to agencies with ≥40 qualifying episodes]
- Panel members would like to have seen data element validation for variables included in the riskadjustment model (and any other critical data elements).

#### <u>Validity</u>

- Validity testing was conducted at the measure score level. The developer also described various data element validation assessments; however, results of these assessments were only summarized, not presented.
- Developers conducted a construct [convergent] validation analysis by correlating (using the Spearman's rank correlation coefficient) the results of this measure with 4 other OASIS performance measures (improvement in ambulation/locomotion, bathing, and bed transfer, and management of oral medications) and the Quality of Patient Care Star Rating measure [it is unclear whether this was a modified version of the measure that excluded the pain measure]
  - o Developers expected statistically significant, strong, positive correlations.
  - $\circ$   $\,$  Correlations with the 4 OASIS measures ranged from 0.51-0.69.
  - Correlation with the star-rating measure = 0.65.
  - These results aligned with supported the developers' hypothesis.
- This measure is risk-adjusted using logistic regression with 114 risk factors (based on 2016 data).
  - Developers discussed previous research linking dual-eligibility status and rural location with use of home health services. They therefore conducted analyses to examine associations between payment source (as a proxy for dual-eligibility) and rurality with this measure. They do include payment source in the risk-adjustment approach, but not rurality.
  - Model discrimination:
    - Overall development sample: c-statistic=0.656
    - Overall model validation sample: c-statistic= 0.657
  - Developers assessed risk-model calibration by calculating McFadden's R<sup>2</sup> and developing risk-decile plots.
    - Overall development sample: McFadden's R<sup>2</sup>=0.053
    - Overall model validation sample: McFadden's R<sup>2</sup>=0.051
- Panel members expressed some concern with excluding transferred patients, questioning whether those patients might have poorer outcomes on this measure. They had a similar concern with excluding patients who died.

## Standing Committee Action Item(s):

• The Standing Committee can discuss reliability and/or validity, or accept the Scientific Methods Panel ratings.

# Questions for the Committee regarding reliability:

• Do you have any concerns that the measure can be consistently implemented (i.e., are measure specifications adequate)?

• The Scientific Methods Panel is satisfied with the reliability testing for the measure. Does the Committee think there is a need to discuss and/or vote on reliability?

#### Questions for the Committee regarding validity:

- Do you have any concerns regarding the validity of the measure (e.g., exclusions, risk-adjustment approach, etc.)?
- The Scientific Methods Panel is satisfied with the validity analyses for the measure. Does the Committee think there is a need to discuss and/or vote on validity?

Preliminary rating for reliability: Preliminary rating for validity:	□ High □ High	⊠ Moderate ⊠ Moderate	□ Low □ Low	<ul><li>Insufficient</li><li>Insufficient</li></ul>			
Scientific Acceptability Prelimina	ry Analysis						
Measure Number: 0177 Measure Title: Improvement in Pai	n Interfering	with Activity					
Type of measure:							
🗆 Process 🛛 Process: Appropri	ate Use 🗆	Structure	Efficiency	Cost/Resource Use			
□⊠ Outcome ⊠□ Outcome: P	RO-PM	Outcome: Inter	mediate Cl	inical Outcome 🛛 Composite			
Data Source:							
Claims      I      Electronic Health	Data 🗆 E	electronic Health	Records	🗆 Management Data			
🗆 🛛 Assessment Data 🛛 🗆 Paper	Medical Rec	ords 🛛 🖓 🖓	trument-Ba	ased Data 🛛 🗆 Registry Data			
Enrollment Data     Other							
Level of Analysis:							
Clinician: Group/Practice     Cl	inician: Indiv	vidual 🛛 🖾 Faci	lity 🗆 He	alth Plan			
Population: Community, County or City Population: Regional and State							
□ Integrated Delivery System □	Other						

#### Measure is:

□ **New** ⊠ **Previously endorsed (**NOTE: Empirical validity testing is expected at time of maintenance review; if not possible, justification is required.)

#### **RELIABILITY: SPECIFICATIONS**

1. Are submitted specifications precise, unambiguous, and complete so that they can be consistently implemented? 
Yes 
No

Submission document: "MIF\_xxxx" document, items S.1-S.22

**NOTE**: NQF staff will conduct a separate, more technical, check of eCQM specifications, value sets, logic, and feasibility, so no need to consider these in your evaluation.

The OASIS items are well specified in the NQF evaluation form.

#### 2. Briefly summarize any concerns about the measure specifications.

Answers are somewhat ambiguous especially as it relates to movement. Data dictionary does not offer any guidance. Measure could also have the unintended consequence of promoting the use of narcotics to score better especially chronic pain patients. One exclusion criterion does cause my concern. Episodes of care ended in transfer to in patient facility were excluded from the denominator because no assessment information was available for these patients. However, it is quite possible that many of these excluded patients might have poor outcomes. Given that a substantial proportion of episodes of care, about 27% (2b2.2), were excluded due to this reason, and particularly if there is across HHAs variation on this exclusion, the measure score may be potentially biased.

None

No concerns.

#### **RELIABILITY: TESTING**

**Submission document:** "MIF\_xxxx" document for specifications, testing attachment questions 1.1-1.4 and section 2a2

- 3. Reliability testing level 🛛 🖾 Measure score 🖾 Data element 🗔 Neither
- 4. Reliability testing was conducted with the data source and level of analysis indicated for this measure ⊠ Yes □ No
- 5. If score-level and/or data element reliability testing was NOT conducted or if the methods used were NOT appropriate, was **empirical <u>VALIDITY</u> testing** of <u>patient-level data</u> conducted?

🗆 Yes 🛛 No

Since data source testing was conducted, appears this question is not answered.

#### 6. Assess the method(s) used for reliability testing

**Submission document:** Testing attachment, section 2a2.2**Submission document:** Testing attachment, section 2a2.2

The Kappa statistic was used to score the inter-rater reliability between the scores at first assessment for admitted or discharged patients. An independent trained team of RNs or physical therapist conducted a separate visit within 24 hours to independently assess the patient. 213 home visits were assessed across 4 states using the entire OASIS survey. For this measure, the underlying questions related to OASIS-C2 item for M1242 (frequency of pain interfering with patient's activity or movement)

The methods used seemed adequate.

For reliability of the performance measure score, the developer tested both measure reliability (test – retest) and facility score reliability (beta-binomial). For the beta-binomial testing, however, it is not clear whether this testing was based on observed results or risk adjusted results. Because this measure is specified as a risk adjusted measure, the testing should be based on risk adjusted results.

<u>SS:</u> The measure developer used appropriate methods to compute reliability estimates at the agency level using a STN model and conducted item-level reliability analyses using kappa agreement to evaluate interrater reliability of assessment items.

JB: Measure score: "...fit a beta-binomial model to estimate measure reliability..." [p8]

"... test-retest reliability using the ICC to measure between-agency variation and within-agency variation..." [p9]

<u>Data element</u>: "...field test of new and existing OASIS items on 12 HHAs in four states for 213 home health patients. Home health registered nurses and physical therapists, trained by the study team, collected data during home visits at start of care (SOC) or resumption of care (ROC), and/or at discharge. Follow-up visits were conducted within 24 hours of the initial field test visit, by a different registered nurse or physical therapist to test interrater reliability...." [p9]

#### 7. Assess the results of reliability testing

Submission document: Testing attachment, section 2a2.3

JB: <u>Measure score</u>: "...AA-ICC is 0.9, and the CA-ICC is also 0.9.

"Beta Binomial Reliability Scores: Mean: 0.95, Median: 0.97 = above 'acceptable' range" [p 9]

<u>Data element</u>: "...The inter-rater reliability (weighted kappa) values for M1242 (frequency of pain interfering with patient's activity or movement) was 0.53 at SOC/ROC and 0.45 at discharge." [p10]

JK: The mean and median inter-rater reliability scores of 0.95 and 0.97 for the entire OASIS survey, were above the range considered acceptable (0.70 - 0.80). Scores for the ambulation section were substantially, inter-rater reliability indicated moderate agreement at SOC/ROC (0.53) and EOC (0.45)

Reliability seemed adequate - details below.

The summary of facility reliability scores (22a.3) showed that 90<sup>th</sup> percentile of facility score reliability is 1.0, indicating that 10% facility scores had reliability of 1.0, this is extremely rare for a risk adjusted measure. It is important to know if these results were based on unadjusted rates.

Weighted kappa is moderate, it would be helpful to report the proportions of agreement as well.

Results indicate both the assessment items and the agency-level scores are sufficiently reliable.

While the raw outcome variable was shown to be reasonably reliable, I have 2 concerns. First, it was not clear whether the risk-adjusted score was modeled in the STN analysis. Given that the measure uses substantial risk adjustment, this seems like a significant omission. Further, item-level reliability was not reported for the 114 variables used in the risk –adjustment model, which also seems like important information to omit.

8. Was the method described and appropriate for assessing the proportion of variability due to real differences among measured entities? NOTE: If multiple methods used, at least one must be appropriate.

Submission document: Testing attachment, section 2a2.2

**⊠Yes ⊠Yes** (ICC method is appropriate, beta-binomial approach is also appropriate but clarification needed for actual testing.) **⊠Yes** 

□No

□Not applicable (score-level testing was not performed)

9. Was the method described and appropriate for assessing the reliability of ALL critical data elements?

Submission document: Testing attachment, section 2a2.2

⊠□Yes

□□□□⊠No- results for risk adjustment variables were not included in item level or score-level testing

□Not applicable (data element testing was not performed)

Denominator exclusions in part include M1700, M1710, M1720. Ideally would have tested these data elements as well. Questionable whether these are 'critical' data elements, e.g. # / % of cases excluded by these OASIS questions.

10. OVERALL RATING OF RELIABILITY (taking into account precision of specifications and <u>all</u> testing results):

⊠ **High** (NOTE: Can be HIGH <u>only if</u> score-level testing has been conducted) (ICC is sufficient for this rating)

⊠ **Moderate** (NOTE: Moderate is the highest eligible rating if score-level testing has <u>not</u> been conducted)

**Low** (NOTE: Should rate <u>LOW</u> if you believe specifications are NOT precise, unambiguous, and complete or if testing methods/results are not adequate)

□⊠Insufficient (NOTE: Should rate <u>INSUFFICIENT</u> if you believe you do not have the information you need to make a rating decision)

# 11. Briefly explain rationale for the rating of OVERALL RATING OF RELIABILITY and any concerns you may have with the approach to demonstrating reliability. Statistical

JK: Statistical testing was appropriate and thorough. Although I thought the lower rate of agreement being much lower than the overall entire form agreement was somewhat concerning.

JB: Sufficient testing performed for measure score & data element. Test results were good. Not able to rate 'high' due to response to #7 above: "...The inter-rater reliability (weighted kappa) values for M1242 (frequency of pain interfering with patient's activity or movement) was 0.53 at SOC/ROC and 0.45 at discharge." Result is only modest.

Only potential concern noted re question #9 above: Denominator exclusions in part include M1700, M1710, M1720. Ideally would have tested these data elements as well. Questionable whether these are 'critical' data elements, e.g. # / % of cases excluded by these OASIS questions.

ZL: Although clarification is needed for the beta-binomial test results and weighted kappa is moderate, test – retest results do indicate high reliability.

SS: Would expect the STN analysis to include the risk adjusted score for each facility and the inter-rater reliability analysis to include the risk adjustment variables as well. Without this information, we cannot tell if the computed (i.e., risk adjusted) scores are reliable.

DN: Methods and results seemed to support reliability of the measure at both data element and measure score levels.

#### VALIDITY: ASSESSMENT OF THREATS TO VALIDITY

#### 12. Please describe any concerns you have with measure exclusions.

Submission document: Testing attachment, section 2b2.

#### No Concerns

One exclusion criterion does cause my concern. Episodes of care ended in transfer to in patient facility were excluded from the denominator because no assessment information was available for these patients. However, it is quite possible that many of these excluded patients might have poor outcomes. Given that a substantial proportion of episodes of care, about 27% (2b2.2), were excluded due to this reason, and particularly if there is across HHAs variation on this exclusion, the measure score may be potentially biased.

Per the completed NQF endorsement form – S.8. Denominator Exclusions: "... the episode of care ended in transfer to inpatient facility or death at home...". In measure specifications we want to avoid excluding cases that may reflect poor quality care. Of course, the quality of care is precisely what we're trying to measure. The concern is a portion of such cases excluded (noted above) may be due to poor quality. Thus, the entity essentially gets a pass on these cases.

None – measure exclusions seemed appropriate and reasonable.

None

# 13. Please describe any concerns you have regarding the ability to identify meaningful differences in performance.

Submission document: Testing attachment, section 2b4.

None

No concern given the percentile distribution on p. 24

None – the developers did a thorough job of addressing the meaningful differences in performance issue.

No

No concerns.

14. Please describe any concerns you have regarding comparability of results if multiple data sources or methods are specified.

Submission document: Testing attachment, section 2b5.

N/A

No concern as it states the only data source is OASIS. It appears OASIS captures everything required for the measure, which includes the exclusions of transfer to inpatient facility or death at home.

N/A

#### 15. Please describe any concerns you have regarding missing data.

Submission document: Testing attachment, section 2b6.

As noted by submitters; ; missing data is rejected and requires resubmission

Given the response to this question it appears there is no missing data: "There are minimal issues with missing data because the OASIS submission system rejects assessments with missing values. The provider must then resubmit the assessment." Thus, no concerns.

None – missing data clearly make a difference, and this will be particularly true for agencies with a small number of episodes from which to generate a score.

No

SS: The measure developer assures us there are 'minimal issues with missing data' as the system apparently rejects forms with missing data. Still, actual rates of missing data would be helpful.

#### 16. Risk Adjustment

- 16a. Risk-adjustment method 🛛 None 🛛 Statistical model 🖓 Stratification
- 16b. If not risk-adjusted, is this supported by either a conceptual rationale or empirical analyses?

□ Yes □ No ⊠□ Not applicable

#### 16c. Social risk adjustment:

16c.1 Are social risk factors included in risk model? □⊠ Yes ⊠□ No □ Not applicable (Payment source as proxy for Medicaid coverage)

16c.2 Conceptual rationale for social risk factors included?

16c.3 Is there a conceptual relationship between potential social risk factor variables and the measure focus? ⊠□ Yes □⊠ No

#### 16d.Risk adjustment summary:

16d.1 All of the risk-adjustment variables present at the start of care? 🛛 Yes 🛛 🗋 No

- 16d.2 If factors not present at the start of care, do you agree with the rationale provided for inclusion? ⊠□ Yes □ No
- 16d.3 Is the risk adjustment approach appropriately developed and assessed? oxtimes Yes  $\hfill\square$  No
- 16d.4 Do analyses indicate acceptable results (e.g., acceptable discrimination and calibration)

🗆 🛛 Yes 🛛 No

16d.5.Appropriate risk-adjustment strategy included in the measure? 
Yes 
No

#### 16e. Assess the risk-adjustment approach Submitters

The only concern I have is the intentional omission of race/ethnicity as a potential adjustment variable. This is a delicate issue, as inclusion of race/ethnicity in a model could create the appearance of having different performance standards for agencies serving different groups, or of excusing poor performance for one or more groups. On the other hand, there have been extensive studies of a relationship between pain experience and pain reporting by members of different racial/ethnic groups going back at least to the 1950s. Race and ethnicity can have an effect on pain reporting that relates to culture more so that any other element of social disadvantage for which race/ethnicity could be a proxy. For this specific measure, it would have been appropriate to formally test race/ethnicity as a potential adjustment variable.

"...The overall model development sample c-statistic is **0.656**. The overall model validation sample cstatistic is **0.657...** The overall model development sample R<sup>2</sup> is **0.053**. The overall model validation sample R<sup>2</sup> is **0.051...**. The plot below shows that the predicted and observed values are similar and monotonically increasing with predicted probability, both of which indicate a well calibrated model. Additionally, we consider the R<sup>2</sup> statistics (included justification for the limited SES in response to 2b3.6) to be sufficient indicators of model fit.." [p22, figure: 22]

JK: Submitters included justification for the limited SES factors used in the model. Logistic regression model was used appropriately to measure the effect of the chosen elements and included only the statistically significant variables in the risk adjusted model.

SS: Overall, the model appears to perform well. However, with 114 risk factors, over-fitting is a real possibility even with this large a sample, but the developer does not address this concern (i.e., use of predicted R-sq).

The adjusted rate is set to 100% if the calculated rate is higher than 100%, the adjusted rate is set to 0% if the calculated rate is lower than 0%.

The developer should articulate the rationale for this approach and report how many facilities were impacted by this approach.

#### For cost/resource use measures ONLY:

- 17. Are the specifications in alignment with the stated measure intent?
  - □ Yes □ Somewhat □ No (If "Somewhat" or "No", please explain)
- 18. Describe any concerns of threats to validity related to attribution, the costing approach, carve outs, or truncation (approach to outliers):

#### **VALIDITY: TESTING**

- 19. Validity testing level: 🛛 Measure score 🔄 Data element 🔤 Both
- 20. Method of establishing validity of the measure score:
  - $\boxtimes \Box$  Face validity
  - **Empirical validity testing of the measure score**
  - □ N/A (score-level testing not conducted)
- 21. Assess the method(s) for establishing validity

#### Submission document: Testing attachment, section 2b2.2 [JB: should refer to 2b1.2]

Spearman rank correlation and expert validity

<u>Measure score</u>: "Convergent validity refers to the extent to which measures that are designed to assess the same construct are related to each other. To evaluate the convergent validity of the measure, Abt calculated the Spearman rank correlations of the *Improvement in Pain Interfering with Activity* measure with other relevant measures, including the publicly-reported measures of home health quality derived from OASIS assessments." [p10]

"...reports the Spearman rank correlation of the *Improvement in Pain Interfering with Activity* measure with a version of the Quality of Patient Care Star Rating, where *Improvement in Pain Interfering with Activity* is excluded from the calculation of the star rating in order to avoid mechanical correlations." [p11]

Methods were generally reasonable, including examining correlations between this measure and other accepted quality of care measures for home health agencies. Since there would be no objective "gold standard" for a patient-reported measure of pain interference, validity testing at the data element level is challenging, but the data reported here support validity, and similar measures of the same concept generally have demonstrated validity.

<u>Data element</u>: Re OASIS: "updated and improved based on input from clinicians and technical experts", "published in the Federal Register for comment... and no objections or suggestions for revision have been noted ..." [p11] "Validity testing included:

1) Consensus validity by expert researcher/clinical panels for outcome measurement and risk factor measurement

2) Consensus validity by expert clinical panels for patient assessment and care planning

3) Criterion or convergent/predictive validity for outcome measurement/risk factor measurement

4) Convergent/predictive validity: case mix adjustment for payment

5) Validation by patient assessment and care planning" [p12]

Spearman rank correlations with other ADL measures were computed for the measure score. However, it is unclear if the correlations use risk-adjusted scores.

Inter-rater reliability scores (weighted kappa) were used as a proxy to establish item-level validity scores for the OASIS item using different raters at 2 different points in time (paired assessments across raters were done within 24 hours). Conceptually, this approach is problematic. For validity testing, one expects a gold standard against which to compare – which assessment is the gold standard in this scenario? Consequently, I find this approach to evaluating item-level validity not compelling. However, score level validity results are appropriately computed.

ZL: Convergent/predictive validity analyses as outlined are reasonable.

#### 22. Assess the results(s) for establishing validity

## Submission document: Testing attachment, section 2b2.3 [JB: should refer to 2b1.3]

Submitters demonstrated strong correlation using the Spearman rank correlation and also expert/ clinical panel

Measure score: 'Spearman rank: 0.50 - 0.69"[p13]

#### Data element:

- 1) Consensus validity: "recommended for measuring patient outcomes..."
- 2) Consensus validity: "recommended for inclusion..."
- 3) Criterion or convergent/predictive validity: "found to be related to other indicators of health status and patient outcomes..."

JB: Note the topic heading here is "data element". The response is not in regard to data element level.

4) Convergent/predictive validity: "Case mix adjustment for payment: The item was tested and is used in the grouping algorithm that, in part, determines the per-episode payment to home health agencies for care provided under the Medicare home health benefit."

JB: Note I don't think case mix adjustment for payment equates to case mix adjustment for risk of an outcome measure. The 3M APR analogy: 1 grouping for severity of illness as it relates to resource consumption, 1 grouping for risk of mortality.

5) Validation by patient assessment and care planning: "reported by practicing clinicians to be effective and useful..." [p13]

Validity results seem adequate

Overall, the results show that the raw measure score is valid given positive correlations with other similar measures of ADL function.

Testing results are acceptable

23. Was the method described and appropriate for assessing conceptually and theoretically sound hypothesized relationships?

Submission document: Testing attachment, section 2b1.

⊠Yes

□No

□Not applicable (score-level testing was not performed)

24. Was the method described and appropriate for assessing the accuracy of ALL critical data elements? *NOTE that data element validation from the literature is acceptable.* 

Submission document: Testing attachment, section 2b1.

⊠□Yes

□⊠No

**Not applicable** (data element testing was not performed)

Item level reliability findings (kappa agreement) do not provide information about the validity of all the data elements required to compute this measure.

# 25. OVERALL RATING OF VALIDITY taking into account the results and scope of all testing and analysis of potential threats.

□ High (NOTE: Can be HIGH only if score-level testing has been conducted)

□ **Moderate** (NOTE: Moderate is the highest eligible rating if score-level testing has NOT been conducted)

**Low** (NOTE: Should rate LOW if you believe that there <u>are</u> threats to validity and/or relevant threats to validity were <u>not assessed OR</u> if testing methods/results are not adequate)

□Insufficient (NOTE: For instrument-based measures and some composite measures, testing at both the score level and the data element level <u>is required</u>; if not conducted, should rate as INSUFFICIENT.)

# 26. Briefly explain rationale for rating of OVERALL RATING OF VALIDITY and any concerns you may have with the developers' approach to demonstrating validity.

I did not have any concerns of the methods used. As noted in questions 21 & 22, the submitters demonstrated solid analysis for validity testing

There are no strong concerns about validity, but a rating of "high" validity would seem to require some more direct evidence of some defined quality of care process measures having a causal relationship with the outcome measure. If the measure indicates that some home health agencies are 'better" than others, what exactly is it that they are doing that is "better"? Then, to what extent does the outcome measure faithfully reflect those differences? A measure with "high" validity would have to be able to demonstrate something like what fraction of the observed variance in the measure score is associated with underlying differences in quality of care, and show that that fraction is large and significant.

This is a valid measure, the main concern I have is with the exclusion criterion that I mentioned earlier.

In general, performed well in testing. Noted as medium vs high due to:

[1] See response to #22: Convergent/predictive validity: "Case mix adjustment for payment: The item was tested and is used in the grouping algorithm that, in part, determines the per-episode payment to home health agencies for care provided under the Medicare home health benefit."

JB: Note I don't think case mix adjustment for payment equates to case mix adjustment for risk of an outcome measure. The 3M APR analogy: 1 grouping for severity of illness as it relates to resource consumption, 1 grouping for risk of mortality.

[2] See response to #22: JB: Note in response to #21 above, CMS notes they used convergent validity. However the results are not noted here.

[3] See response to #22 – specifically #3 under "data element" heading: JB: Note the topic heading here is "data element". The response is not in regard to data element level.

Although not clear if risk adjusted score was used to determine measure score correlations and the approach to determining item-level validity was not sufficient, the outcome variable's correlations with other ADLs (presumably raw score) indicates this measure has sufficient validity at the score-level.

#### FOR COMPOSITE MEASURES ONLY: Empirical analyses to support composite construction

- 27. What is the level of certainty or confidence that the empirical analysis demonstrates that the component measures add value to the composite and that the aggregation and weighting rules are consistent with the quality construct?
  - □High

□Moderate

- Low
- □Insufficient
- 28. Briefly explain rationale for rating of EMPIRICAL ANALYSES TO SUPPORT COMPOSITE CONSTRUCTION

## ADDITIONAL RECOMMENDATIONS

29. If you have listed any concerns in this form, do you believe these concerns warrant further discussion by the multi-stakeholder Standing Committee? If so, please list those concerns below.

## **Committee Pre-evaluation Comments:**

# Criteria 2: Scientific Acceptability of Measure Properties (including all 2a, 2b, and 2c)

## 2a1. Reliability-Specifications

Comments:

- No concerns
- I appreciate the Scientific Methods Panel's evaluation and the staff's recommendation of moderate reliability and validity.
- They present kappa values of 0.45 at the start of care and 0.53 at discharge which are considered moderate. I wonder what explains the difference between start of care and discharge. The mean signal-to-noise reliability estimate is 0.95. The split sample estimate is IRR(2,1)= 0.90; IRR(3,1)= 0.90
- No concerns.
- I am happy to accept the evaluation of the Scientific Methods Panel and have nothing to add.
- The description of the measure (The percentage of home health episodes of care during which the frequency of the patient's pain when moving around improved.) talks about the relationship between pain and the ability to move around. I am not sure how movement is being captured by the measure.

## 2a2. Reliability-Testing

Comments:

- No concerns
- No
- Curious why the difference in kappa values between start of care and discharge
- No concerns. Recommendation of the Scientific Advisory Panel can be accepted

- No
- The reliability test results appear to be good and so I have no concerns other than that noted in question 5.

### 2b1. Validity-Testing

#### Comments:

- No concerns, but would like to know more about the specific exclusions and how they may have impacted the results and why they were excluded.
- No
- Construct convergent validity was done by comparing results with 4 other OASIS performance measures (improvement in ambulation/locomotion, bathing, and bed transfer, and management of oral medications). These may or may not be closely related to pain. correlation ranged from 0.51 0.69.
- No concerns.
- No
- I have no concerns.

#### 2b4-7. Threats to Validity

Comments:

- same answer as above in 7.
- My only concern is the existence of psychosocial, emotional and psychoogical variables that may interfere with pain management. I'm not sure if this measure would get to them.
- Only demonstrated improvements in pain. Does not consider those without or with minimal pain at the start. Meaningful differences are noted by substantial interquartile range. Also, measure performance is related to other metrics in the direction expected with statistically significant differences in measure performance across strata. Minimal issues with missing data because OASIS submission system rejects assessments with missing values.
- No apparent threats to validity in these categories.
- I am happy to accept the evaluation of the Scientific Methods Panel and have nothing to add.
- The analyses conducted provide evidence that the measure appears to identify meaningful differences. I
  base this on the numbers that reveal that there is still room for improvement, the TCV reveals variation in
  measure distribution, the interquartile range indicates the presence of measure spread, and the same
  information for HHAs stratifed by census region. Only one set of specification so no discussion inregard to
  qustion 2b5. Missing data not an issue due to the OASIS submission system requirement that no
  assessments are accepted if they have missing data.

#### 2b2-3. Other Threats to Validity

Comments:

- yes conceptual relationship shown
- Threats to validity include unknown addiction, presense drug seekers in the home, comorbid conditions that generate pain.
- I worry about excluding patient without or with minimal pain at the start. There is a risk adjust strategy presented
- As with the other measures in this group, some concern about the exclusion of patients who were hospitalized as a result of this encounter. This is a substantial group in this care. Possibly the developers can sample this group and report data on outcomes.
- I have no concerns regarding validity aside from those mentioned by the Scientific Methods Panel, regarding the exclusion of transferred and deceased patients. I do not believe that this exclusion constitutes a threat to the validity of the measure.
- I believe thought should be given to the exclusion of episodes ending in transfer to inpatient facilities. Oftentimes, patients are transferred when their pain is out of control and there is need for inpatient care to address it. I have no concerns about risk-adjustment.

# Criterion 3. Feasibiilty

#### Maintenance measures - no change in emphasis - implementation issues may be more prominent

**<u>3. Feasibility</u>** is the extent to which the specifications including measure logic, require data that are readily available or could be captured without undue burden and can be implemented for performance measurement.

• The data for this measure comes from the OASIS dataset. OASIS captures assessment information during the home health episode of care. Collection and transmission of OASIS is a requirement for the Medicare Home Health Conditions of Participation.

#### **Questions for the Committee:**

- Are the required data elements routinely generated and used during care delivery?
- Are the required data elements available in electronic form, e.g., EHR or other electronic sources?

Preliminary rating for feasibility: 🛛 High 🗌 Moderate 🗌 Low 🔲 Insufficient

#### **Committee Pre-evaluation Comments: Criteria 3: Feasibility**

#### 3. Feasibility

Comments:

- no concerns
- Measurement during a home health encounter is entirely appropriate.
- Data come from the OASIS dataset. Collecting this data is a requiriement of the Medicare Home health conditions of Participation. I am not aware of the extent to which there is compliance with this requirement.
- All data elements are routinely generated and available in electronic form.
- I have no concerns regarding data collection. It appears to be easily captured without undue burden.
- Other than the one concern I have about movement, I do not see any issues with feasibility.

## Criterion 4: Usability and Use

<u>Maintenance measures</u> – increased emphasis – much greater focus on measure use and usefulness, including both impact/improvement and unintended consequences

#### 4a. Use (4a1. Accountability and Transparency; 4a2. Feedback on measure)

<u>4a. Use</u> evaluate the extent to which audiences (e.g., consumers, purchasers, providers, policymakers) use or could use performance results for both accountability and performance improvement activities.

**4a.1.** Accountability and Transparency. Performance results are used in at least one accountability application within three years after initial endorsement and are publicly reported within six years after initial endorsement (or the data on performance results are available). If not in use at the time of initial endorsement, then a credible plan for implementation within the specified timeframes is provided.

#### Current uses of the measure

Publicly reported?	🛛 Yes 🛛	Νο
Current use in an accountability program?	🛛 Yes 🛛	No 🗆 UNCLEAR
OR		
Planned use in an accountability program?	🗆 Yes 🗆	No

#### Accountability program details

- The measure is used in the following:
  - Home Health Compare (public reporting)
  - Home Health Star Ratings (internal quality improvement). Agencies receive a "Outcome Quality Measure Report" that allows agencies to benchmark their performance against other agencies across the state and nationally, as well as their own performance from prior time periods.
  - It is not clear from the submission whether this measure is also included in the Home Health Quality Reporting Program (HHQRP) and Home Health Value-Based Purchasing (HHVBP) program.

**4a.2. Feedback on the measure by those being measured or others.** Three criteria demonstrate feedback: 1) those being measured have been given performance results or data, as well as assistance with interpreting the measure results and data; 2) those being measured and other users have been given an opportunity to provide feedback on the measure performance or implementation; 3) this feedback has been considered when changes are incorporated into the measure

#### Feedback on the measure by those being measured or others

- The developer <u>reports</u> that home health agencies obtain feedback on the measure via quarterly Quality of Patient Care Star Rating Provider Preview Reports. Agencies are able to review for errors or submit questions via email. Additionally, HHQRP training was conducted for agencies in 2017.
- While the developer did not summarize the feedback from home health agencies, they did note that no requests for modifications have been made.

#### Additional Feedback:

• No additional feedback has been provided.

#### Questions for the Committee:

- How have (or can) the performance results be used to further the goal of high-quality, efficient healthcare?
- How has the measure been vetted in real-world settings by those being measured or others?

#### Preliminary rating for Use: 🛛 Pass 🗌 No Pass

#### 4b. Usability (4a1. Improvement; 4a2. Benefits of measure)

<u>**4b.**</u> <u>Usability</u> evaluate the extent to which audiences (e.g., consumers, purchasers, providers, policymakers) use or could use performance results for both accountability and performance improvement activities.

**4b.1 Improvement.** Progress toward achieving the goal of high-quality, efficient healthcare for individuals or populations is demonstrated.

#### Improvement results

• The developer provided data to show improvement results. The developer <u>also described</u> improvement over time within population subgroups. They noted that the large improvements from from 2015 to 2016 likely were due to the introduction of several initiatives that incorporate this measure: the Quality of Patient Care (QoPC) Star Ratings, a composite of this measure and several others that has been publicly reported on Home Health Compare since July 2015 and Home Health Value Based Purchasing (HHVBP).

**4b2. Benefits vs. harms.** Benefits of the performance measure in facilitating progress toward achieving highquality, efficient healthcare for individuals or populations outweigh evidence of unintended negative consequences to individuals or populations (if such evidence exists).

#### Unexpected findings (positive or negative) during implementation

• The developer stated that recent improvement in this measure has been relatively large compared to historical trends due to the implementation of two initiatives that involve this measure (the QoPC Star Ratings and HHVBP, beginning in 2015 and 2016, respectively).

#### **Potential harms**

• The developer did not indicate any potential harms or benefits from this measure.

#### **Questions for the Committee:**

• Do the benefits of the measure outweigh any potential unintended consequences?

Preliminary rating for Usability and use:	🛛 High	🛛 Moderate	🗆 Low	Insufficient
---	--------	------------	-------	--------------

## Committee Pre-evaluation Comments: Criteria 4: Usability and Use

#### 4a1. Use-Accountability and Transparency

#### Comments:

- Agree that the measure is in use and provides feedback and comparison for national and regional benchmarking and consumer info on quality. Not sure on opportunity for feedback from HHA. Would be helpful to have that information.
- No requested modifications.
- Public reported with Home Health Compare. Agencies can benchmark performance against other agencies across the state and nationally. Agencies can give feedback via quarterly quality of patient care star rating provider preview reports. The result were not summarized, but they did note that no request for modifications have been made
- Users have been given feedback and opportunity to provide feedback. Developers report that no changes have thus far been suggested.
- The measure is used in Home Health Compare, for public reporting, Home Health Star Ratings for internal quality improvement allowing agencies to compare their performance to one another. Although the developer did not summarize the feedback they have obtained feedback and noted that no requests for modifications have been made.
- Measures are publicly reported via the CMS Home Health Compare website and their Outcome Quality Measure Report. HHAs can submit feedback via an email box provided for that purpose. Thus far, no modification requests have been made.

#### 4b2. Usability-Improvement

Comments:

- Yes, the benefits outweigh any potential harms
- Public reporting and accountability are benefits.
- Data are provided that show improvement results and over time within population subgroups. Benefits are to facilitate better care. Harms--misses deterioration in pain
- Measure is in use as a performance improvement measure. No harms apparent.
- Large improvements from 2015 to 2016 were reported by the developers most likely due to several initiatives that incorporate the measure. The developer did not indicate any potential harms or benefits from this measure.
- Developer provided evidence of improvement overall and for certain subgroups. I do not see evidence of any unintended consequences.

## Criterion 5: Related and Competing Measures

#### **Related measures**

• 0209: Comfortable Dying: Pain Brought to a Comfortable Level Within 48 Hours of Initial Assessment [facility-level outcome measure in ambulatory (hospice) setting]

#### Harmonization

• NQF may ask the Committee to make recommendations for combining or harmonizing measures.

### **Committee Pre-evaluation Comments: Criterion 5: Related and Competing Measures**

#### 5. Related and Competing Measures

Comments:

- there is one related measure NQF#0209 would like to discuss whether harmonization is appropriate even though that is more specific and the exclusions may be different.
- I don't find Comfortable Ding (0209) to be a competing measure. People who are learning pain management in relation to activity are in a different disease state than those who are in hospice care.
- The developer did not report related measure but the pre-evaluation committee noted 0.09: comfortable dying: Pain brought to a comfortable level withing 48 hours of initial assessment. No harmonization reported
- No related or competing measures apparent.
- There is one related measure 0209 on Comfortable Dying, it does not appear to compete with this measure.
- There are no related or competing measures and so harmonization is not necessary.

# **Public and Member Comments**

No NQF members have submitted support/non-support choices as of January 25, 2019. No comments have been submitted as of January 25, 2019.

Additional evaluations and submission materials attachments...

## **1.** Evidence and Performance Gap – Importance to Measure and Report

Extent to which the specific measure focus is evidence-based, important to making significant gains in healthcare quality, and improving health outcomes for a specific high-priority (high-impact) aspect of healthcare where there is variation in or overall less-than-optimal performance. *Measures must be judged to meet all sub criteria to pass this criterion and be evaluated against the remaining criteria.* 

#### 1a. Evidence to Support the Measure Focus – See attached Evidence Submission Form

#### nqf\_evidence\_attachment\_7.1--PAIN-jsr.docx

# 1a.1 <u>For Maintenance of Endorsement:</u> Is there new evidence about the measure since the last update/submission?

Do not remove any existing information. If there have been any changes to evidence, the Committee will consider the new evidence. Please use the most current version of the evidence attachment (v7.1). Please use red font to indicate updated evidence.

Yes

#### 1a. Evidence (subcriterion 1a)

Measure Number (if previously endorsed): 0177

Measure Title: Improvement in pain interfering with activity

# IF the measure is a component in a composite performance measure, provide the title of the Composite Measure here:

#### Date of Submission:

Instructions

- Complete 1a.1 and 1a.2 for all measures. If instrument-based measure, complete 1a.3.
- Complete EITHER 1a.2, 1a.3 or 1a.4 as applicable for the type of measure and evidence.
- For composite performance measures:
  - A separate evidence form is required for each component measure unless several components were studied together.
  - If a component measure is submitted as an individual performance measure, attach the evidence form to the individual measure submission.
- All information needed to demonstrate meeting the evidence subcriterion (1a) must be in this form. An appendix of *supplemental* materials may be submitted, but there is no guarantee it will be reviewed.
- If you are unable to check a box, please highlight or shade the box for your response.
- Contact NQF staff regarding questions. Check for resources at <u>Submitting Standards webpage</u>.

<u>Note</u>: The information provided in this form is intended to aid the Standing Committee and other stakeholders in understanding to what degree the evidence for this measure meets NQF's evaluation criteria.

1a. Evidence to Support the Measure Focus

The measure focus is evidence-based, demonstrated as follows:

• <u>Outcome</u>: <u>a</u> Empirical data demonstrate a relationship between the outcome and at least one healthcare structure, process, intervention, or service. If not available, wide variation in performance can be used as

evidence, assuming the data are from a robust number of providers and results are not subject to systematic bias.

- <u>Intermediate clinical outcome</u>: a systematic assessment and grading of the quantity, quality, and consistency of the body of evidence <u>b</u> that the measured intermediate clinical outcome leads to a desired health outcome.
- <u>Process</u>: <u>c</u> a systematic assessment and grading of the quantity, quality, and consistency of the body of evidence <u>b</u> that the measured process leads to a desired health outcome.
- <u>Structure</u>: a systematic assessment and grading of the quantity, quality, and consistency of the body of evidence <u>b</u> that the measured structure leads to a desired health outcome.
- <u>Efficiency</u>: <u>d</u> evidence not required for the resource use component.
- For measures derived from <u>patient reports</u>, evidence should demonstrate that the target population values the measured outcome, process, or structure and finds it meaningful.
- <u>Process measures incorporating Appropriate Use Criteria</u>: See NQF's guidance for evidence for measures, in general; guidance for measures specifically based on clinical practice guidelines apply as well. Notes

a. Generally, rare event outcomes do not provide adequate information for improvement or discrimination; however, serious reportable events that are compared to zero are appropriate outcomes for public reporting and quality improvement.

b. The preferred systems for grading the evidence are the Grading of Recommendations, Assessment, Development and Evaluation (GRADE) guidelines and/or modified GRADE.

c. Clinical care processes typically include multiple steps: assess  $\rightarrow$  identify problem/potential problem  $\rightarrow$  choose/plan intervention (with patient input)  $\rightarrow$  provide intervention  $\rightarrow$  evaluate impact on health status. If the measure focus is one step in such a multistep process, the step with the strongest evidence for the link to the desired outcome should be selected as the focus of measurement. Note: A measure focused only on collecting PROM data is not a PRO-PM.

d. Measures of efficiency combine the concepts of resource use <u>and</u> quality (see NQF's <u>Measurement</u> <u>Framework: Evaluating Efficiency Across Episodes of Care</u>; <u>AQA Principles of Efficiency Measures</u>).

**1a.1.This is a measure of**: (should be consistent with type of measure entered in De.1)

Outcome

⊠ Outcome:

□Patient-reported outcome (PRO):

PROs include HRQoL/functional status, symptom/symptom burden, experience with care, healthrelated behaviors. (A PRO-based performance measure is not a survey instrument. Data may be collected using a survey instrument to construct a PRO measure.)

□ Intermediate clinical outcome (*e.g., lab value*):

 $\Box$  Process:

- □ Appropriate use measure:
- $\Box$  Structure:
- $\Box$  Composite:
- **1a.2 LOGIC MODEL** Diagram or briefly describe the steps between the healthcare structures and processes (e.g., interventions, or services) and the patient's health outcome(s). The relationships in the diagram should be easily understood by general, non-technical audiences. Indicate the structure, process or outcome being measured.

Appropriate home health care interventions should improve the rates of improvement in frequency of pain when moving around.

**1a.3 Value and Meaningfulness: IF** this measure is derived from patient report, provide evidence that the target population values the measured *outcome, process, or structure* and finds it meaningful. (Describe how and from whom their input was obtained.)

#### Not applicable

\*\*RESPOND TO ONLY ONE SECTION BELOW -EITHER 1a.2, 1a.3 or 1a.4) \*\*

1a.2 FOR OUTCOME MEASURES including PATIENT REPORTED OUTCOMES - Provide empirical data demonstrating the relationship between the outcome (or PRO) to at least one healthcare structure, process, intervention, or service.

Pain, both acute and chronic, has Both acute and chronic pain have been identified as an area requiring health care provider intervention. Many patients who receive home health care experience pain, which can interfere with activity and affect virtually all aspects of a patient's daily life, as well as impact other aspects of health status. Appropriate evaluation and management of pain is recognized as very important to the well-being of patients. Clinical practice guidelines identify effective interventions for chronic pain including both pharmacologic and non-pharmacologic. High-quality home health care appropriately evaluates and manages pain, and reduction in pain can be considered a marker of high-value care for patients with pain. This measure meets National Priorities Partnership (NPP) Goals related to access to effective treatment for relief of suffering from symptoms such as pain.

Pain is a serious problem with adverse impact on a wide range of outcomes from functional capacity, to quality of life and mortality.<sup>1 2 3</sup> Pain, both acute and chronic, is prevalent among community-dwelling older adults. Gaps in assessment and management of pain persist.<sup>4 5</sup>

<sup>&</sup>lt;sup>1</sup> Cornelius R, Herr KA, Gordon DB, Kretzer K, Butcher HK. Evidence-Based Practice Guideline: Acute Pain Management in Older Adults. J Gerontol Nurs. 2017 Feb 1;43(2):18-27.

<sup>&</sup>lt;sup>2</sup> Crowe M, Gillon D, Jordan J, McCall C. Older peoples' strategies for coping with chronic non-malignant pain: A qualitative meta-synthesis. Int J Nurs Stud. 2017 Mar;68:40-50.

<sup>&</sup>lt;sup>3</sup> Kradag AS, Bakan AB, Varol E, Aslan G. Investigation of pain and life satisfaction in older adults. Geriatr Gerontol Int. 2017 Jul 28.

<sup>&</sup>lt;sup>4</sup> Henchoz Y, Büla C, Guessous I, Rodondi N, Goy R, Demont M, & Santos-Eggimann B. (2017). Chronic symptoms in a representative sample of community-dwelling older people: a cross-sectional study in Switzerland. BMJ Open, 7(1), e014485.

<sup>&</sup>lt;sup>5</sup> Nawai A, Leveille SG, Shmerling RH, van der Leeuw G, Bean JF. Pain severity and pharmacologic pain management among community-living older adults: the MOBILIZE Boston study. Aging Clin Exp Res. 2017 Dec;29(6):1139-1147.

Experience of bothersome pain may interfere with functional improvement during rehabilitation, and meeting rehabilitation goals for community-dwelling older adults.<sup>6</sup> Pain is associated with higher risk of falls<sup>7 8 9</sup> and activity-limiting pain is associated with recurrent falls.<sup>10</sup> Additional adverse impacts of pain reported in research include more rapid cognitive decline<sup>11</sup> and depressive symptoms.<sup>12 13</sup> Evidence about the association between pain and emergency department (ED) use is mixed. In one study of homeless adults, severe pain was identified as one factor associated with higher emergency department (ED) use, <sup>14</sup> although in another pain was not reported at significantly different frequency among community-dwelling older adults who used the ED versus those who did not.<sup>15</sup>

Multiple conditions are predictors or consequences of pain. Arthritis and depression are significantly associated with longstanding pain that substantially limits participation in daily living<sup>16</sup> Increasing frequency of falls, fatigue and having depression are predictors for presence of severe daily pain among community-dwelling older adults.<sup>17</sup>

<sup>8</sup> Talarska D, Strugała M, Szewczyczak M, Tobis S, Michalak M, Wróblewska I, Wieczorowska-Tobis K. Is independence of older adults safe considering the risk of falls? BMC Geriatr. 2017 Mar 14;17(1):66.

<sup>9</sup> Kendall JC, Hvid LG, Hartvigsen J, Fazalbhoy A, Azari MF, Skjødt M, Robinson SR, Caserotti P. Impact of musculoskeletal pain on balance and concerns of falling in mobility-limited, community-dwelling Danes over 75 years of age: a cross-sectional study. Aging Clin Exp Res. 2017 Dec 11.

<sup>10</sup> Agudelo-Botero M, Giraldo-Rodríguez L, Murillo-González JC, Mino-León D, Cruz-Arenas E. Factors associated with occasional and recurrent falls in Mexican community-dwelling older people. PLoS One. 2018 Feb 20;13(2):e0192926. doi: 10.1371/journal.pone.0192926. eCollection 2018.

<sup>11</sup> Whitlock EL, Diaz-Ramirez LG, Glymour MM, Boscardin WJ, Covinsky KE, Smith AK. Association Between Persistent Pain and Memory Decline and Dementia in a Longitudinal Cohort of Elders. JAMA Intern Med. 2017 Jun 5.

<sup>12</sup> Xiang X, Brooks J. Correlates of depressive symptoms among homebound and semi-homebound older adults. J Gerontol Soc Work. 2017 Jan 27.

<sup>13</sup> Sugai K, Takeda-Imai F, Michikawa T, Nakamura T, Takebayashi T, Nishiwaki Association Between Knee Pain, Impaired Function, and Development of Depressive Symptoms. J Am Geriatr Soc. 2018 Mar;66(3):570-576. doi: 10.1111/jgs.15259. Epub 2018 Feb 14.

<sup>14</sup> Raven MC, Tieu L, Lee CT, Ponath C, Guzman D, Kushel M. Emergency Department Use in a Cohort of Older Homeless Adults: Results From the HOPE HOME Study. Acad Emerg Med. 2017 Jan;24(1):63-74.

<sup>15</sup> Dermody G, Sawyer P, Kennedy R, Williams C, Brown CJ. ED Utilization and Self-Reported Symptoms in Community-Dwelling Older Adults. J Emerg Nurs. 2017 Jan;43(1):57-69.

<sup>16</sup> Janevic MR, McLaughlin SJ, Heapy AA, Thacker C, Piette JD. Racial and Socioeconomic Disparities in Disabling Chronic Pain: Findings from the Health and Retirement Study. J Pain. 2017 Jul 28.

<sup>17</sup> Crowe M, Jordan J, Gillon D, McCall C, Frampton C, Jamieson H. The prevalence of pain and its relationship to falls, fatigue, and depression in a cohort of older people living in the community. J Adv Nurs. 2017 May 5.

<sup>&</sup>lt;sup>6</sup> Gell NM, Mroz TM, Patel KV. Rehabilitation Services Use and Patient Reported Outcomes among Older Adults in the United States. Arch Phys Med Rehabil. 2017 Apr 3.

<sup>&</sup>lt;sup>7</sup> Kitayuguhi J, Kamada M, Inoue S, Kamioka H, Abe T, Okada S, Mutoh Y. Association of low back and knee pain with falls in Japanese community-dwelling older adults: A 3-year prospective cohort study. Geriatr Gerontol Int. 2017 Jun;17(6):875-884.

Some recent evidence suggests racial disparities in the pain experience, with African American older adults experiencing unmet pain needs<sup>18</sup> and reporting more pain-related disability<sup>19</sup> than other racial/ethnic groups. Murtaugh and colleagues<sup>20</sup> found that Hispanics and non-Hispanic blacks reported a greater number of pain sites, worse pain intensity and higher levels of pain-related disability than non-Hispanic whites and others. This cross-sectional study of 588 patients with activity-limiting pain, who were admitted to home care for physical therapy, also provided evidence that race/ethnicity interacts with pain self-efficacy and depressive symptoms in their association with mean pain intensity and pain-related disability, respectively.

Non-pharmacological interventions, as well as pharmacologic management, may have a positive effect on pain management in older adults.<sup>21</sup> For example, an 8-week chair yoga program provided to older adults with lower extremity osteoarthritis was associated with reduction in pain interference, and this effect was sustained 3 months post-intervention.<sup>22</sup> Exercise training combined with psychosocial intervention may be more effective than exercise training alone for management of chronic pain in community-dwelling older adults.<sup>23</sup> Evidence suggests that self-management assessment and support should be tailored by pain condition.<sup>24</sup>

1a.3. SYSTEMATIC REVIEW(SR) OF THE EVIDENCE (for INTERMEDIATE OUTCOME, PROCESS, OR STRUCTURE PERFORMANCE MEASURES, INCLUDING THOSE THAT ARE INSTRUMENT-BASED) If the evidence is not based on a systematic review go to section 1a.4) If you wish to include more than one systematic review, add additional tables.

What is the source of the <u>systematic review of the body of evidence</u> that supports the performance measure? A systematic review is a scientific investigation that focuses on a specific question and uses explicit, prespecified scientific methods to identify, select, assess, and summarize the findings of similar but separate studies. It may include a quantitative synthesis (meta-analysis), depending on the available data. (IOM)

□ Clinical Practice Guideline recommendation (with evidence review)

□ US Preventive Services Task Force Recommendation

□ Other systematic review and grading of the body of evidence (*e.g., Cochrane Collaboration, AHRQ Evidence Practice Center*)

 $\Box$  Other

<sup>20</sup> Murtaugh CM, Beissner KL, Barrón Y, Trachtenberg MA, Bach E, Henderson CR Jr, Sridharan S, Reid MC. Pain and Function in Home Care: A Need for Treatment Tailoring to Reduce Disparities? Clin J Pain. 2017 Apr;33(4):300-309.

<sup>21</sup>Horgas A. Pain management in older adults. Nur Clin N Am 52 (2017) eq-e7. http://dx.doi.org/10/1016/j.cnur.2017.08.001

<sup>22</sup> Park J, McCaffrey R, Newman D, Liehr P, Ouslander JG. A Pilot Randomized Controlled Trial of the Effects of Chair Yoga on Pain and Physical Function Among Community-Dwelling Older Adults With Lower Extremity Osteoarthritis. J Am Geriatr Soc. 2017 Mar;65(3):592-597.

<sup>23</sup> Hirase T, Kataoka H, Nakano J, Inokuchi S, Sakamoto J, Okita M. Impact of frailty on chronic pain, activities of daily living and physical activity in community-dwelling older adults: A cross-sectional study. Geriatr Gerontol Int. 2018 Mar 26.

<sup>24</sup> Mann EG, Harrison MB, LeFort S, VanDenKerkhof EG. A Canadian Survey of Self-Management Strategies and Satisfaction with Ability to Control Pain: Comparison of Community Dwelling Adults with Neuropathic Pain versus Adults with Non-neuropathic Chronic Pain. Pain Manag Nurs. 2018 Mar 1. pii: S1524-9042(16)30242-9. doi: 10.1016/j.pmn.2017.10.016.

<sup>&</sup>lt;sup>18</sup> Robinson-Lane Robinson-Lane SG, Vallerand AH. Pain Treatment Practices of Community-Dwelling Black Older Adults. Pain Manag Nurs. 2017 Dec 13. pii: S1524-9042(17)30422-8.

<sup>&</sup>lt;sup>19</sup> Janevic MR, McLaughlin SJ, Heapy AA, Thacker C, Piette JD. Racial and Socioeconomic Disparities in Disabling Chronic Pain: Findings from the Health and Retirement Study. J Pain. 2017 Jul 28.

Source of Systematic Review:	
• Title	
Author	
Date	
Citation, including page number	
• URL	
Quote the guideline or recommendation verbatim about the process, structure or intermediate outcome being measured. If not a guideline, summarize the conclusions from the SR.	
Grade assigned to the <b>evidence</b> associated with the recommendation with the definition of the grade	
Provide all other grades and definitions from the evidence grading system	
Grade assigned to the <b>recommendation</b> with definition of the grade	
Provide all other grades and definitions from the recommendation grading system	
Body of evidence:	
<ul> <li>Quantity – how many studies?</li> </ul>	
<ul> <li>Quality – what type of studies?</li> </ul>	
Estimates of benefit and consistency across studies	
What harms were identified?	
Identify any new studies conducted since the SR. Do the new studies change the conclusions from the SR?	

#### **1a.4 OTHER SOURCE OF EVIDENCE**

If source of evidence is NOT from a clinical practice guideline, USPSTF, or systematic review, please describe the evidence on which you are basing the performance measure.

Four studies were found to support the need to measure pain in an adult home health care population.

**1a.4.1 Briefly SYNTHESIZE the evidence that supports the measure.** A list of references without a summary is not acceptable.

#### 1a.4.2 What process was used to identify the evidence?

PubMed and Google Scholar searches were performed using key word "Home health care" in combination with each of the key word: "Pain," The search was limited to 2006 – present.

#### 1a.4.3. Provide the citation(s) for the evidence.

A: 1) Maxwell, C.J., Dalby, D.M., Slater, M., Patten, S.B., Hogan, D.B., Eliasziw, M. & Hirdes, J.P. (2008). The prevalence and management of current daily pain among older home care patients. *Pain, 138*, 208-216. 2) This study was a cross-sectional analysis of prevalence and correlates of pharmacotherapy for daily pain in 2779 elders (≥ 65) receiving home care in Ontario from 1999 – 2001. C) 47.8% of patients reported daily pain and

21.6% of patients reporting pain received no analgesia. D) Findings from this study indicate that pain affects a large proportion of patients receiving home health care services, and that pain assessment and intervention is critical for this population.

B). Ornstein, K., Wajnberg, A., Kaye-Kauderer, H., Winkel, G., DeCherrie, L., Zhang, M. & Soriano, T. (2013). Reduction in symptoms for homebound patients receiving home-based primary and palliative care. *Journal of Palliative Medicine, 16*(9), 1048-1054. 2)This study evaluated the effectiveness of a home-based palliative care program consisting of a home visit, physician follow up, social work, and in-home visits as needed. The study sample consisted of 140 homebound patients and data were collected using the Edmonton Symptom Assessment Scale at 3 and 12 weeks following program enrollment via telephone. 3) Pain and other symptoms were reduced at 3 weeks, with reductions maintained at 12 weeks. 4) This study highlights the role that home care interventions can play in reduction in symptom burden, particularly for pain.

C: 1) Bach, E., Beissner, K., Murtaugh, C., Trachtenberg, M., & Carrington Reid (M.) (2013). Implementing a cognitive-behavioral pain self-management in home health care part 2: Feasibility and acceptability study. *J Geriatr Phys Ther, 36*(3), 130-137. 2) The study tested a new physical therapy cognitive-behavioral intervention for pain management in sample of 21 home health care patients. Patients were interviewed following each intervention session and asked whether they found the techniques to be effective in relieving pain. 3) Patient ratings of techniques learned as part of the intervention ranged from 71.4 - 81.2% in terms of helpfulness for pain management. 4) While the sample was small and there was no control group, this study suggests that home care interventions may be helpful in helping patients to manage pain. Thus, measurement of pain is important for a home health care population.

D: 1) Zyczkowska, J., Szczerbinska, K., Jantzi, M.R. & Hirdes, J.P. (2007). Pain among the oldest old in community and institutional settings. *Pain, 129,* 167-176. 2) This cross-sectional study used evaluated pain ratings in 193,158 home health care and complex continuing care (similar to skilled nursing facility) patients in Ontario, Canada. Patients were grouped by age into 5-year categories (e.g., 65-69 yo, etc.). The sample included 788 patients 100 years old or older. Reports of pain were highest for the youngest-old group (65-69; mean rating of 1.41, 95% CI = 1.38 - 1.43)), declining in each age group, with the lowest mean occurring in the 100 - 115 yo group (mean rating of .98; 95% CI = .88 - 1.07). In home care, arthritis and osteoporosis diagnoses were associated with highest pain scores. For home care patients, the odds ratio for pain was .76 (95% CI = .73 - .80) for the 70-74 age group, lowering to was 0.42 (CL: .35 - .51) in the oldest group. 4) This study suggests that pain is a problem with an older adult home care population, which should be addressed through careful assessment and evaluation of pain management interventions.

#### 1b. Performance Gap

Demonstration of quality problems and opportunity for improvement, i.e., data demonstrating:

- considerable variation, or overall less-than-optimal performance, in the quality of care across providers; and/or
- Disparities in care across population groups.

**1b.1. Briefly explain the rationale for this measure** (*e.g.*, how the measure will improve the quality of care, the benefits or improvements in quality envisioned by use of this measure)

<u>If a COMPOSITE</u> (e.g., combination of component measure scores, all-or-none, any-or-none), SKIP this question and answer the composite questions.

Both acute and chronic pain have been identified as areas requiring health care provider intervention. Many patients who receive home health care experience pain, which can interfere with activity and affect virtually all aspects of a patient's daily life, as well as

impact other aspects of health status. Appropriate evaluation and management of pain is recognized as very important to the wellbeing of patients. Clinical practice guidelines identify effective interventions for chronic pain including both pharmacologic and nonpharmacologic.

High-quality home health care appropriately evaluates and manages pain, and reduction in pain can be considered a marker of high-value care for patients with pain.

**1b.2.** Provide performance scores on the measure as specified (<u>current and over time</u>) at the specified level of analysis. (<u>This is required for maintenance of endorsement</u>. Include mean, std dev, min, max, interquartile range, scores by decile. Describe the data source including number of measured entities; number of patients; dates of data; if a sample, characteristics of the entities include.) This information also will be used to address the sub-criterion on improvement (4b1) under Usability and Use.

See attachment, "Importance to Report" for a tabular presentation of these data, as well as the tables below.

Tables 1 and 2 in the "Importance to Report" attachment show observed and predicted measure performance, respectively, for calendar years 2010 through 2016, including the number of HHAs and the average number of episodes for HHAs. For each table, the top panel shows this information for all HHAs with at least one episode for which the measure is available. The bottom panel shows this information for HHAs with at least 20 episode for which the measure is available.

Calendar	Number	Average	ННА	Std.	Minimum	10th	25th	50th	75th	90th	Maximum	IQR*
Year	of HHAs	Episodes	Average	Dev.		Percentile	Percentile	Percentile	Percentile	Percentile		
		per HHA										
HHAs wi	th >=1 Va	lid Episod	de									
2010	10,831	236	61.0%	21.0%	0.0%	33.3%	51.4%	63.3%	73.2%	85.7%	100.0%	21.8%
2011	11,487	278	61.9%	20.6%	0.0%	35.0%	51.8%	63.8%	74.4%	86.9%	100.0%	22.6%
2012	11,748	273	61.9%	21.4%	0.0%	33.3%	51.4%	64.3%	75.0%	88.0%	100.0%	23.6%
2013	11,893	282	62.2%	21.9%	0.0%	33.3%	51.0%	64.4%	76.1%	89.2%	100.0%	25.1%
2014	11,832	300	61.6%	22.3%	0.0%	31.0%	50.3%	64.2%	75.9%	88.9%	100.0%	25.6%
2015	11,527	335	62.6%	22.9%	0.0%	30.3%	50.7%	66.5%	77.5%	90.0%	100.0%	26.8%
2016	11,166	347	65.6%	23.7%	0.0%	31.6%	54.5%	70.0%	82.1%	92.9%	100.0%	27.6%
HHAs wi	th >=20 \	/alid Episc	ode									
2010	8,504	299	63.8%	15.7%	0.0%	44.4%	55.6%	64.6%	73.0%	83.4%	100.0%	17.4%
2011	9,476	335	63.8%	16.8%	0.0%	42.9%	55.0%	64.7%	74.0%	85.0%	100.0%	19.0%
2012	9,664	330	64.3%	17.7%	0.0%	42.1%	55.2%	65.4%	75.2%	86.6%	100.0%	20.0%
2013	9,749	342	64.5%	18.2%	0.0%	41.0%	55.0%	65.5%	76.0%	87.5%	100.0%	21.0%
2014	9,609	367	64.2%	18.4%	0.0%	40.0%	54.6%	65.4%	75.9%	87.3%	100.0%	21.3%
2015	9,462	406	65.6%	19.0%	0.0%	40.7%	55.6%	67.8%	77.8%	88.8%	100.0%	22.2%
2016	9,028	427	69.4%	18.9%	0.0%	44.8%	60.0%	71.8%	82.6%	91.7%	100.0%	22.6%

Table 1.Observed HHA-level Performance on Improvement in Pain Interfering with Activity by Calendar Year

\*The IQR (interquartile range) is a measure of variability. It is calculated by subtracting the 25th percentile value from the 75th percentile value.

Table 2. Risk Adjusted HHA-level Performance on Improvement in Pain Interfering with Activity by Calendar Year

Calendar	Number	Average	нна	Std.	Minimum	10th	25th	50th	75th	90th	Maximum	IQR*
Year	of HHAs	Episodes	Average	Dev.		Percentile	Percentile	Percentile	Percentile	Percentile		
		per HHA										
HHAs wit	th >=1 Va	lid Episod	le					-				-
2010	10,831	236	61.9%	19.6%	0.0%	36.3%	53.4%	64.2%	73.5%	84.6%	100.0%	20.1%
2011	11,487	278	63.0%	19.1%	0.0%	38.3%	54.0%	64.8%	74.7%	85.8%	100.0%	20.7%
2012	11,748	273	63.2%	19.8%	0.0%	37.2%	53.9%	65.4%	75.7%	87.1%	100.0%	21.8%
2013	11,893	282	63.6%	20.2%	0.0%	36.1%	53.8%	65.8%	76.6%	88.5%	100.0%	22.8%
2014	11,832	300	63.3%	20.7%	0.0%	34.7%	53.3%	65.4%	76.5%	88.6%	100.0%	23.2%
2015	11,527	335	64.4%	21.2%	0.0%	34.4%	54.0%	67.5%	78.2%	90.2%	100.0%	24.1%
2016	11,166	347	67.7%	21.6%	0.0%	36.2%	57.7%	71.2%	82.5%	93.2%	100.0%	24.8%
HHAs wi	th >=20 V	/alid Episc	ode									
2010	8,504	299	64.4%	14.8%	0.0%	46.4%	56.9%	65.1%	73.1%	82.3%	100.0%	16.2%
2011	9,476	335	64.6%	15.8%	0.0%	44.9%	56.7%	65.5%	74.3%	83.9%	100.0%	17.6%
2012	9,664	330	65.1%	16.5%	0.0%	44.9%	56.8%	66.2%	75.4%	85.8%	100.0%	18.6%
2013	9,749	342	65.5%	17.0%	0.0%	44.1%	56.7%	66.6%	76.2%	86.9%	100.0%	19.5%
2014	9,609	367	65.3%	17.3%	0.0%	43.3%	56.6%	66.4%	76.2%	86.7%	100.0%	19.7%
2015	9,462	406	66.8%	17.9%	0.0%	43.7%	57.8%	68.6%	78.1%	88.8%	100.0%	20.4%
2016	9,028	427	70.7%	17.5%	0.0%	48.2%	62.1%	72.6%	82.6%	91.9%	100.0%	20.6%

\*The IQR (interquartile range) is a measure of variability. It is calculated by subtracting the 25th percentile value from the 75th percentile value.

Table 3 provides characteristics of all home health patients in 2016 for which this measure could be calculated.

Table 3. Patients Characteristics - All Patients in Measure Calculation, 2016

Population Group		# of Patients	% of Patients	
Total		3,876,352	100.0%	
Gender	Male	1,412,799	36.4%	
	Female	2,463,553	63.6%	
Race	White	2,980,418	76.9%	
	Black	504,404	13.0%	
	Hispanic	285,073	7.4%	
	Other	106,457	2.7%	
Age	Under 65	691,217	17.8%	
	65-74	1,087,413	28.1%	
	75-84	1,149,717	29.7%	
	85 and Over	948,005	24.5%	
Disability Status	No	3,012,247	77.7%	
	Yes	864,105	22.3%	
Dual Enrollment in	No	2,965,087	76.5%	
Medicare and Medicaid	Yes	911,265	23.5%	

1b.3. If no or limited performance data on the measure as specified is reported in 1b2, then provide a summary of data from the literature that indicates opportunity for improvement or overall less than optimal performance on the specific focus of measurement.

#### See attachment, "Importance to Report" for a tabular presentation of these data.

**1b.4.** Provide disparities data from the measure as specified (current and over time) by population group, e.g., by race/ethnicity, gender, age, insurance status, socioeconomic status, and/or disability. (*This is required for maintenance of endorsement*. Describe the data source including number of measured entities; number of patients; dates of data; if a sample, characteristics of the entities included.) For measures that show high levels of performance, i.e., "topped out", disparities data may demonstrate an opportunity for improvement/gap in care for certain sub-populations. This information also will be used to address the sub-criterion on improvement (4b1) under Usability and Use.

Tables 4 and 5 in the "Importance to Report" attachment, and below, show observed and predicted measure performance for population groups, respectively. Measure performance improved from 2013 to 2016 for all population groups. For some population groups, performance gaps between subgroups also diminished over time.

For example, for gender the difference between observed measure performance for males and females in 2013 was 1.9 percentage points. This difference slightly decreased to 1.4 percentage points in 2016. The difference in measure performance between those in the northeast versus those in the south also decreased over time.

For some population groups, disparities did increase. For example, the difference in measure performance between patients who were not dually eligible for Medicaid compared to those who were increased from 1.8 in 2013 percentage points to 3.1 percentage points in 2016. The difference in performance between small and large agencies also widened over time.

<b>Population Group</b>	o	2013	2014	2015	2016
All Episodes		67.5%	67.7%	70.1%	74.4%
Gender	Male	68.8%	68.9%	71.1%	75.3%
	Female	66.9%	67.1%	69.5%	74.0%
Race	White	67.1%	67.4%	70.2%	74.8%
	Black	66.7%	66.9%	69.0%	73.1%
	Hispanic	72.2%	71.1%	70.8%	73.4%
	Other	70.1%	70.5%	72.1%	75.3%
Age	Under 65	62.5%	62.8%	65.1%	69.5%
	65-74	68.5%	68.6%	70.9%	75.3%
	75-84	68.9%	69.1%	71.5%	75.8%
	85 and Over	68.6%	68.8%	71.3%	75.5%
Disability Status	No	69.0%	69.2%	71.5%	75.6%
	Yes	63.2%	63.3%	65.6%	70.3%
Dual Enrollment	No	68.0%	68.3%	70.9%	75.2%
in Medicare and Medicaid	Yes	66.2%	66.2%	67.9%	72.0%
Agency Size	Small	55.0%	54.0%	53.1%	56.2%
	Medium	67.6%	67.2%	68.5%	71.7%
	Large	67.7%	68.1%	70.7%	75.3%
Census Region	Northeast	68.6%	69.0%	71.5%	75.3%
	Midwest	68.4%	68.6%	70.2%	73.4%
	South	66.2%	66.3%	69.0%	74.5%
	West	67.9%	68.1%	70.5%	74.4%

Table 4. Observed Episode-Level Measure Performance by Population Group

Table 5. Predicted Episode-Level Measure Performance by Population Group

Population Group		2013	2014	2015	2016
All Episodes		66.9%	67.1%	65.5%	67.3%
Gender	Male	67.0%	67.1%	65.5%	67.1%
	Female	66.9%	67.0%	65.5%	67.4%
Race	White	67.2%	67.4%	65.9%	67.8%
	Black	65.3%	65.4%	63.6%	65.2%
	Hispanic	67.1%	66.8%	64.6%	66.1%
	Other	67.7%	67.6%	65.4%	67.1%
Age	Under 65	62.7%	62.9%	61.3%	62.9%
	65-74	69.1%	69.2%	67.8%	69.6%
	75-84	67.7%	67.8%	66.3%	68.1%
	85 and Over	66.7%	66.8%	65.0%	66.8%
Disability Status	No	67.7%	67.7%	66.1%	67.9%
	Yes	64.7%	65.0%	63.5%	65.3%
Dual Enrollment	No	67.4%	67.5%	66.0%	67.8%
in Medicare and Medicaid	Yes	65.7%	65.8%	64.0%	65.7%
Agency Size	Small	64.2%	64.0%	62.2%	63.3%
	Medium	66.4%	66.3%	64.4%	66.3%
	Large	67.1%	67.3%	65.8%	67.6%
Census Region	Northeast	66.9%	67.1%	65.4%	67.1%
	Midwest	67.5%	67.6%	66.0%	67.7%
	South	66.6%	66.8%	65.2%	67.1%
	West	67.1%	67.0%	65.6%	67.4%

1b.5. If no or limited data on disparities from the measure as specified is reported in 1b.4, then provide a summary of data from the literature that addresses disparities in care on the specific focus of measurement. Include citations. Not necessary if performance data provided in 1b.4

See 1.b4

# 2. Reliability and Validity—Scientific Acceptability of Measure Properties

Extent to which the measure, <u>as specified</u>, produces consistent (reliable) and credible (valid) results about the quality of care when implemented. *Measures must be judged to meet the sub criteria for both reliability and validity to pass this criterion and be evaluated against the remaining criteria.* 

**2a.1. Specifications** The measure is well defined and precisely specified so it can be implemented consistently within and across organizations and allows for comparability. eMeasures should be specified in the Health Quality Measures Format (HQMF) and the Quality Data Model (QDM).

**De.5. Subject/Topic Area** (check all the areas that apply):

Palliative Care and End-of-Life Care

**De.6. Non-Condition Specific** (check all the areas that apply):

Health and Functional Status : Change

**De.7. Target Population Category** (Check all the populations for which the measure is specified and tested if any):

#### Elderly, Populations at Risk : Individuals with multiple chronic conditions

**S.1. Measure-specific Web Page** (*Provide a URL link to a web page specific for this measure that contains current detailed specifications including code lists, risk model details, and supplemental materials. Do not enter a URL linking to a home page or to general information.*)

https://www.cms.gov/Medicare/Quality-Initiatives-Patient-Assessment-Instruments/HomeHealthQualityInits/Home-Health-Quality-Measures.html

**S.2a.** <u>If this is an eMeasure</u>, HQMF specifications must be attached. Attach the zipped output from the eMeasure authoring tool (MAT) - if the MAT was not used, contact staff. (Use the specification fields in this online form for the plain-language description of the specifications)

#### This is not an eMeasure Attachment:

**S.2b. Data Dictionary, Code Table, or Value Sets** (and risk model codes and coefficients when applicable) must be attached. (Excel or csv file in the suggested format preferred - if not, contact staff)

Attachment Attachment: isc\_mstr\_-V2.21.1-\_FINAL\_08-15-2017-636776316361945348.xlsx

**S.2c.** Is this an instrument-based measure (i.e., data collected via instruments, surveys, tools, questionnaires, scales, etc.)? Attach copy of instrument if available.

Attachment Attachment: OASIS-C2-AllItems-10-2016-636686582612898016.pdf

**S.2d.** Is this an instrument-based measure (i.e., data collected via instruments, surveys, tools, questionnaires, scales, etc.)? Attach copy of instrument if available.

#### Clinician

**S.3.1.** For maintenance of endorsement: Are there changes to the specifications since the last updates/submission. If yes, update the specifications for S1-2 and S4-22 and explain reasons for the changes in S3.2.

#### No

**S.3.2.** For maintenance of endorsement, please briefly describe any important changes to the measure specifications since last measure update and explain the reasons.

#### No significant changes

**S.4. Numerator Statement** (Brief, narrative description of the measure focus or what is being measured about the target population, i.e., cases from the target population with the target process, condition, event, or outcome) DO NOT include the rationale for the measure.

*IF an OUTCOME MEASURE, state the outcome being measured. Calculation of the risk-adjusted outcome should be described in the calculation algorithm (S.14).* 

The number of home health episodes of care where the value recorded on the discharge assessment indicates less frequent pain at discharge than at start (or resumption) of care.

**S.5. Numerator Details** (All information required to identify and calculate the cases from the target population with the target process, condition, event, or outcome such as definitions, time period for data collection, specific data collection items/responses, code/value sets – Note: lists of individual codes with descriptors that exceed 1 page should be provided in an Excel or csv file in required format at S.2b)

*IF an OUTCOME MEASURE,* describe how the observed outcome is identified/counted. Calculation of the riskadjusted outcome should be described in the calculation algorithm (S.14).

The number of home health episodes where the value recorded for the OASIS-C2 item M1242 ("Frequency of Pain Interfering with Activity") on the discharge assessment is numerically less than the value recorded on the start (or resumption) of care assessment, indicating less frequent pain interfering with activity at discharge.

**S.6. Denominator Statement** (Brief, narrative description of the target population being measured)

Number of home heath episodes of care ending with a discharge during the reporting period, other than those covered by generic or measure- specific exclusions.

**S.7. Denominator Details** (All information required to identify and calculate the target population/denominator such as definitions, time period for data collection, specific data collection items/responses, code/value sets – Note: lists of individual codes with descriptors that exceed 1 page should be provided in an Excel or csv file in required format at S.2b.)

*IF an OUTCOME MEASURE, describe how the target population is identified. Calculation of the risk-adjusted outcome should be described in the calculation algorithm (S.14).* 

All home health episodes of care (except those defined in the denominator exclusions) in which the patient was eligible to improve in pain interfering with activity or movement (i.e., were not at the optimal level of health status according to the "Frequency of Pain Interfering" OASIS-C2 item M1242).

**S.8. Denominator Exclusions** (Brief narrative description of exclusions from the target population)

All home health episodes where there is no pain reported at the start (or resumption) of care assessment, or the patient is non-responsive, or the episode of care ended in transfer to inpatient facility or death at home, or the episodes is covered by one of the generic exclusions.

**S.9. Denominator Exclusion Details** (All information required to identify and calculate exclusions from the denominator such as definitions, time period for data collection, specific data collection items/responses, code/value sets – Note: lists of individual codes with descriptors that exceed 1 page should be provided in an Excel or csv file in required format at S.2b.)

Home health episodes of care for which [1] at start/resumption of care OASIS item M1242 = 0, indicating the patient had no pain; OR [2] at start/ resumption of care, OASIS item M1700 "Cognitive Functioning" is 4, or M1710 "When Confused" is NA, or M1720 "When Anxious" is NA, indicating the patient is non-responsive; OR [3] The patient did not have a discharge assessment because the episode of care ended in transfer to inpatient facility or death at home; OR [4] All episodes covered by the generic exclusions:

- a. Pediatric home health patients less than 18 years of age as data are not collected for these patients.
- b. Home health patients receiving maternity care only.
- c. Home health clients receiving non-skilled care only.
- d. Home health patients for which neither Medicare nor Medicaid are a payment source.
- e. The episode of care does not end during the reporting period.
- f. If the agency sample includes fewer than 20 episodes after all other patient-level exclusions are applied, or if the agency has been in operation less than six months, then the data is suppressed from public reporting on Home Health Compare.

**S.10. Stratification Information** (Provide all information required to stratify the measure results, if necessary, including the stratification variables, definitions, specific data collection items/responses, code/value sets, and the risk-model covariates and coefficients for the clinically-adjusted version of the measure when appropriate – Note: lists of individual codes with descriptors that exceed 1 page should be provided in an Excel or csv file in required format with at S.2b.)

#### Not Applicable

**S.11. Risk Adjustment Type** (Select type. Provide specifications for risk stratification in measure testing attachment)

Statistical risk model

If other:

S.12. Type of score:

Rate/proportion

If other:

**S.13. Interpretation of Score** (*Classifies interpretation of score according to whether better quality is associated with a higher score, a lower score, a score falling within a defined interval, or a passing score*)

#### Better quality = Higher score

**S.14. Calculation Algorithm/Measure Logic** (*Diagram or describe the calculation of the measure score as an ordered sequence of steps including identifying the target population; exclusions; cases meeting the target process, condition, event, or outcome; time period for data, aggregating data; risk adjustment; etc.*)

1. Define an episode of care (the unit of analysis): Data from matched pairs of OASIS assessments for each episode of care (start or resumption of care paired with a discharge or transfer to inpatient facility) are used to calculate individual patient outcome measures.

2. Identify target population: All episodes of care ending during a specified time interval (usually a period of twelve months), subject to generic and measure-specific exclusions.

Generic exclusions: Episodes of care ending in discharge due to death (M0100\_ASSMT\_REASON[2] = 08).

Measure specific exclusions: Episodes of care ending in transfer to inpatient facility (M0100\_ASSMT\_REASON[2] IN (06,07), patients who are comatose or non-responsive at start/resumption of care (M1700\_COG\_FUNCTION[1] = 04 OR M1710\_WHEN\_CONFUSED[1] = NA OR M1720\_WHEN\_ANXIOUS[1] = NA), and patients with no pain interfering with activity at start/resumption of care (M1242\_PAIN\_FREQ\_ACTVTY\_MVMT [1] = 00).

Cases meeting the target outcome are those where the patient has less pain interfering with activity at discharge than at start/resumption of care:

M1242\_PAIN\_FREQ\_ACTVTY\_MVMT[2] < M1242\_PAIN\_FREQ\_ACTVTY\_MVMT[1].

3. Aggregate the Data: The observed outcome measure value for each HHA is calculated as the percentage of cases meeting the target population (denominator) criteria that meet the target outcome (numerator) criteria.

4. Risk Adjustment: The expected probability for a patient is calculated using the following formula:

P(x)=1/(1+e^(-(a+?¦?b\_i x\_i ?)))

Where:

P(x) = predicted probability of achieving outcome x

a = constant parameter listed in the model documentation

bi = coefficient for risk factor i in the model documentation

xi = value of risk factor i for this patient. See the attached zipped risk adjustment file for detailed lists and specifications of risk factors.

Predicted probabilities for all patients included in the measure denominator are then averaged to derive an expected outcome value for the agency. This expected value is then used, together with the observed (unadjusted) outcome value and the expected value for the national population of home health agency patients for the same data collection period, to calculate a risk-adjusted outcome value for the home health agency. The formula for the adjusted value of the outcome measure is as follows:

X(A\_ra )= X(A\_obs )+ X(N\_exp )-X(A\_exp)

Where:

X(Ara) = Agency risk-adjusted outcome measure value

X(Aobs) = Agency observed outcome measure value

X(Aexp) = Agency expected outcome measure value

X(Nexp) = National expected outcome measure value

If the result of this calculation is a value greater than 100%, the adjusted value is set to 100%. Similarly, if the result is a negative number the adjusted value is set to zero.

**S.15. Sampling** (*If measure is based on a sample, provide instructions for obtaining the sample and guidance on minimum sample size.*)

<u>IF an instrument-based</u> performance measure (e.g., PRO-PM), identify whether (and how) proxy responses are allowed.

#### Not Applicable

**S.16. Survey/Patient-reported data** (*If measure is based on a survey or instrument, provide instructions for data collection and guidance on minimum response rate.*)

Specify calculation of response rates to be reported with performance measure results.

#### Not Applicable

**S.17. Data Source** (Check ONLY the sources for which the measure is SPECIFIED AND TESTED).

If other, please describe in S.18.

#### **Electronic Health Data**

**S.18. Data Source or Collection Instrument** (Identify the specific data source/data collection instrument (e.g. name of database, clinical registry, collection instrument, etc., and describe how data are collected.)

<u>IF instrument-based</u>, identify the specific instrument(s) and standard methods, modes, and languages of administration.

The measure is calculated based on the data obtained from the Home Health Outcome and Assessment Information Set (OASIS), which is a statutorily required core standard assessment instrument that home health agencies integrate into their own patient-specific, comprehensive assessment to identify each patient's need for home care. The instrument is used to collect valid and reliable information for patient assessment, care planning, and service delivery in the home health setting, as well as for the home health quality assessment and performance improvement program. Home health agencies are required to collect OASIS data on all nonmaternity Medicare/Medicaid patients, 18 or over, receiving skilled services. Data are collected at specific time points (admission, resumption of care after inpatient stay, recertification every 60 days that the patient remains in care, transfer, death, and at discharge). HH agencies are required to encode and transmit patient OASIS data to the OASIS repositories Each HHA has on-line access to outcome and process measure reports based on their own OASIS data submissions, as well as comparative state and national aggregate reports, case mix reports, and potentially avoidable event reports. CMS regularly collects OASIS data for storage in the national OASIS repository, and makes measures based on these data (including the Improvement in Pain Interfering with Activity measure) available to consumers and to the general public through the Medicare Home Health Compare website.

The current version of OASIS is OASIS C2. Starting January 1, 2019, OASIS D will be in effective. Differences include added, deleted, modified items and responses.

**S.19. Data Source or Collection Instrument** (available at measure-specific Web page URL identified in S.1 OR in attached appendix at A.1)

Available at measure-specific web page URL identified in S.1

**S.20. Level of Analysis** (Check ONLY the levels of analysis for which the measure is SPECIFIED AND TESTED) Facility

**S.21. Care Setting** (Check ONLY the settings for which the measure is SPECIFIED AND TESTED)

#### Home Care

If other:

**S.22.** <u>COMPOSITE Performance Measure</u> - Additional Specifications (Use this section as needed for aggregation and weighting rules, or calculation of individual performance measures if not individually endorsed.)

#### Not Applicable

#### 2. Validity – See attached Measure Testing Submission Form

Testing\_Form\_Pain\_20180730.docx,RiskAdjustmentModel-636686587298525074.zip

#### 2.1 For maintenance of endorsement

Reliability testing: If testing of reliability of the measure score was not presented in prior submission(s), has reliability testing of the measure score been conducted? If yes, please provide results in the Testing attachment. Please use the most current version of the testing attachment (v7.1). Include information on all testing conducted (prior testing as well as any new testing); use red font to indicate updated testing.

#### Yes

#### 2.2 For maintenance of endorsement

Has additional empirical validity testing of the measure score been conducted? If yes, please provide results in the Testing attachment. Please use the most current version of the testing attachment (v7.1). Include information on all testing conducted (prior testing as well as any new testing); use red font to indicate updated testing.

Yes

#### 2.3 For maintenance of endorsement

Risk adjustment: For outcome, resource use, cost, and some process measures, risk-adjustment that includes social risk factors is not prohibited at present. Please update sections 1.8, 2a2, 2b1,2b4.3 and 2b5 in the Testing attachment and S.140 and S.11 in the online submission form. NOTE: These sections must be updated even if social risk factors are not included in the risk-adjustment strategy. You MUST use the most current version of the Testing Attachment (v7.1) -- older versions of the form will not have all required questions.

Yes - Updated information is included

Measure Testing (subcriteria 2a2, 2b1-2b6)

Measure Number (*if previously endorsed*): 0177 Measure Title: Improvement in Pain Interfering with Activity Date of Submission: <u>8/1/2018</u>

#### Type of Measure:

Outcome ( <i>including PRO-PM</i> )	Composite – STOP – use composite testing form
Intermediate Clinical Outcome	Cost/resource
Process (including Appropriate Use)	Efficiency
□ Structure	

Instructions

- Measures must be tested for all the data sources and levels of analyses that are specified. *If there is more than one set of data specifications or more than one level of analysis, contact NQF staff* about how to present all the testing information in one form.
- For <u>all</u> measures, sections 1, 2a2, 2b1, 2b2, and 2b4 must be completed.
- For outcome and resource use measures, section 2b3 also must be completed.

- If specified for <u>multiple data sources/sets of specificaitons</u> (e.g., claims and EHRs), section 2b5 also must be completed.
- Respond to <u>all</u> questions as instructed with answers immediately following the question. All information on testing to demonstrate meeting the subcriteria for reliability (2a2) and validity (2b1-2b6) must be in this form. An appendix for *supplemental* materials may be submitted, but there is no guarantee it will be reviewed.
- If you are unable to check a box, please highlight or shade the box for your response.
- Maximum of 25 pages (*incuding questions/instructions;* minimum font size 11 pt; do not change margins). *Contact NQF staff if more pages are needed.*
- Contact NQF staff regarding questions. Check for resources at <u>Submitting Standards webpage</u>.
- For information on the most updated guidance on how to address social risk factors variables and testing in this form refer to the release notes for version 7.1 of the Measure Testing Attachment.

<u>Note</u>: The information provided in this form is intended to aid the Standing Committee and other stakeholders in understanding to what degree the testing results for this measure meet NQF's evaluation criteria for testing.

2a2. Reliability testing <u>e</u> demonstrates the measure data elements are repeatable, producing the same results a high proportion of the time when assessed in the same population in the same time period and/or that the measure score is precise. For instrument-based measures (including PRO-PMs) and composite performance measures, reliability should be demonstrated for the computed performance score.

2b1. Validity testing <u>f</u> demonstrates that the measure data elements are correct and/or the measure score correctly reflects the quality of care provided, adequately identifying differences in quality. For instrument-based measures (including PRO-PMs) and composite performance measures, validity should be demonstrated for the computed performance score.

2b2. Exclusions are supported by the clinical evidence and are of sufficient frequency to warrant inclusion in the specifications of the measure; g

## AND

If patient preference (e.g., informed decisionmaking) is a basis for exclusion, there must be evidence that the exclusion impacts performance on the measure; in such cases, the measure must be specified so that the information about patient preference and the effect on the measure is transparent (e.g., numerator category computed separately, denominator exclusion category computed separately). <u>h</u>

2b3. For outcome measures and other measures when indicated (e.g., resource use):

• an evidence-based risk-adjustment strategy (e.g., risk models, risk stratification) is specified; is based on patient factors (including clinical and social risk factors) that influence the measured outcome and are present at start of care; <u>ij</u> and has demonstrated adequate discrimination and calibration

OR

• rationale/data support no risk adjustment/ stratification.

2b4. Data analysis of computed measure scores demonstrates that methods for scoring and analysis of the specified measure allow for identification of statistically significant and practically/clinically meaningful differences in performance;

OR

there is evidence of overall less-than-optimal performance.

2b5. If multiple data sources/methods are specified, there is demonstration they produce comparable results.

2b6. Analyses identify the extent and distribution of missing data (or nonresponse) and demonstrate that performance results are not biased due to systematic missing data (or differences between responders and nonresponders) and how the specified handling of missing data minimizes bias.

Notes

e. Reliability testing applies to both the data elements and computed measure score. Examples of reliability testing for data elements include, but are not limited to: inter-rater/abstractor or intra-rater/abstractor studies; internal consistency for multi-item scales; test-retest for survey items. Reliability testing of the measure score addresses precision of measurement (e.g., signal-to-noise).

f. Validity testing applies to both the data elements and computed measure score. Validity testing of data elements typically analyzes agreement with another authoritative source of the same information. Examples of validity testing of the measure score include, but are not limited to: testing hypotheses that the measures scores indicate quality of care, e.g., measure scores are different for groups known to have differences in quality assessed by another valid quality measure or method; correlation of measure scores with another valid indicator of quality for the specific topic; or relationship to conceptually related measures (e.g., scores on process measures to scores on outcome measures). Face validity of the measure score as a quality indicator may be adequate if accomplished through a systematic and transparent process, by identified experts, and explicitly addresses whether performance scores resulting from the measure as specified can be used to distinguish good from poor quality. The degree of consensus and any areas of disagreement must be provided/discussed.

g. Examples of evidence that an exclusion distorts measure results include, but are not limited to: frequency of occurrence, variability of exclusions across providers, and sensitivity analyses with and without the exclusion.

h. Patient preference is not a clinical exception to eligibility and can be influenced by provider interventions.

i. Risk factors that influence outcomes should not be specified as exclusions.

j. With large enough sample sizes, small differences that are statistically significant may or may not be practically or clinically meaningful. The substantive question may be, for example, whether a statistically significant difference of one percentage point in the percentage of patients who received smoking cessation counseling (e.g., 74 percent v. 75 percent) is clinically meaningful; or whether a statistically significant difference of \$25 in cost for an episode of care (e.g., \$5,000 v. \$5,025) is practically meaningful. Measures with overall less-than-optimal performance may not demonstrate much variability across providers.

#### NOTE: ALL TESTING CONDUCTED IN THIS FORM RELY UPON MORE RECENT DATA AND AN UPDATED RISK ADJUSTMENT MODEL COMPARED TO THE PREVIOUS NQF SUBMISSION. WE DO NOT MARK ANY RESPONSES IN RED BECAUSE MOST RESPONSES WERE UPDATED.

## 1. DATA/SAMPLE USED FOR <u>ALL</u> TESTING OF THIS MEASURE

Often the same data are used for all aspects of measure testing. In an effort to eliminate duplication, the first five questions apply to all measure testing. If there are differences by aspect of testing, (e.g., reliability vs. validity) be sure to indicate the specific differences in question 1.7.

**1.1.** What type of data was used for testing? (Check all the sources of data identified in the measure specifications and data used for testing the measure. Testing must be provided for <u>all</u> the sources of data specified and intended for measure implementation. If different data sources are used for the numerator and denominator, indicate N [numerator] or D [denominator] after the checkbox.)

Measure Specified to Use Data From: (must be consistent with data sources entered in S.17)	Measure Tested with Data From:
$\Box$ abstracted from paper record	$\square$ abstracted from paper record
🗆 claims	🗆 claims
□ registry	□ registry
$\Box$ abstracted from electronic health record	$\square$ abstracted from electronic health record
eMeasure (HQMF) implemented in EHRs	eMeasure (HQMF) implemented in EHRs
🖾 other: Electronic Clinical Data	🛛 other: Electronic Clinical Data

**1.2.** If an existing dataset was used, identify the specific dataset (*the dataset used for testing must be consistent with the measure specifications for target population and healthcare entities being measured; <i>e.g., Medicare Part A claims, Medicaid claims, other commercial insurance, nursing home MDS, home health OASIS, clinical registry*).

Home Health OASIS-C2

1.3. What are the dates of the data used in testing? January 1, 2016 to December 31, 2016

**1.4. What levels of analysis were tested**? (testing must be provided for <u>all</u> the levels specified and intended for measure implementation, e.g., individual clinician, hospital, health plan)

Measure Specified to Measure Performance of: (must be consistent with levels entered in item S.20)	Measure Tested at Level of:
🗆 individual clinician	🗆 individual clinician
□ group/practice	□ group/practice
⊠ hospital/facility/agency	⊠ hospital/facility/agency
🗆 health plan	🗆 health plan
🗆 other:	□ other:

# 1.5. How many and which <u>measured entities</u> were included in the testing and analysis (by level of analysis and data source)? (*identify the number and descriptive characteristics of measured entities included in the analysis (e.g., size, location, type); if a sample was used, describe how entities were selected for inclusion in the sample*)

To calculate the intra-class correlation (ICC) as part of reliability testing, the measure developer included Medicare-certified agencies with at least 40 home health quality episodes ending between January 1, 2016 and December 31, 2016 and meeting the measure denominator criteria.<sup>25</sup> There were 7,862 such agencies (70.4 percent of the 11,166 agencies with at least one quality episode meeting the measure denominator criteria ending during the same time period). The sample included all quality episodes at these agencies (3,825,093 in total) meeting the measure denominator criteria ending between January 1, 2016 to December 31, 2016.

To calculate the beta-binomial scores (as part of reliability testing) and conduct analyses related to validity testing and exclusions, the measure developer included Medicare-certified agencies with at least 20 home health quality episodes ending between January 1, 2016 and December 31, 2016 and meeting the measure denominator criteria. There were 9,028 such agencies (80.9 percent of the 11,166 agencies with at least one quality episode meeting the measure denominator criteria ending during the same time period). The sample included all quality episodes meeting the measure denominator criteria at these agencies (3,858,808 in total) ending between January 1, 2016 to December 31, 2016.

The risk adjustment model was developed using OASIS national repository data from assessments submitted between January 1, 2016 and December 31, 2016 (which included ~ 6.4 million episodes of care).

# **1.6.** How many and which <u>patients</u> were included in the testing and analysis (by level of analysis and data source)? (*identify the number and descriptive characteristics of patients included in the analysis (e.g., age, sex, race, diagnosis); if a sample was used, describe how patients were selected for inclusion in the sample)*

The table below identifies the patients by population group used to calculate the intra-class correlation (ICC) as part of reliability testing. As noted in section 1.5, these are the patients represented in Medicare-certified

<sup>&</sup>lt;sup>25</sup> A minimum of 40 episodes is used instead of the 20 episode criteria for public reporting because the ICC requires splitting each HHA into two samples. To ensure that each sample has a 20 episode minimum, we use a 40 episode minimum for the HHA when evaluating test-retest reliability.

agencies with at least 40 home health quality episodes ending between January 1, 2016 and December 31, 2016 and meeting the measure denominator criteria, the data represented 3,015,107 patients.

Population Group		# of Patients	% of Patients
Total		3,015,107	100%
Gender	Male	1,107,978	36.75%
	Female	1,907,129	63.25%
Race	White	2,347,495	77.86%
	Black	372,234	12.34%
	Hispanic	212,524	7.05%
	Other	82,854	2.75%
Age	Under 65	539,105	17.88%
	65-74	875,488	29.04%
	75-84	893,741	29.64%
	85 and Over	706,773	23.44%
Dual Enrollment in Medicare	No	2,373,618	78.72%
and Medicaid	Yes	641,489	21.28%
Currently or Originally Eligible	No	2,379,618	78.92%
for Medicare due to Disability	Yes	635,489	21.08%
Location of HHA by Census	Northeast	629,733	20.88%
Region	Midwest	639,420	21.21%
	South	1,204,152	39.94%
	West	530,012	17.58%
	Missing	11,790	0.39%

Number/Percentage of Patients Represented in HHAs with At Least 40 Valid Episodes, By Population Group

The table below identifies the patients by population group used to calculate the beta-binomial scores (as part of reliability testing) and conduct analyses related to validity testing and exclusions. As noted in section 1.5, these are the patients represented in Medicare-certified agencies with at least 20 home health quality episodes ending between January 1, 2016 and December 31, 2016 and meeting the measure denominator criteria, the data represented 3,052,331 patients.

Population Group		# of Patients	% of Patients
Total		3,052,331	100%
Gender	Male	1,121,526	36.74%
	Female	1,930,805	63.26%
Race	White	2,366,813	77.54%
	Black	381,384	12.49%
	Hispanic	219,305	7.19%
	Other	84,829	2.78%
Age	Under 65	548,023	17.95%
	65-74	885,724	29.02%
	75-84	904,034	29.62%
	85 and Over	714,550	23.41%

#### Number/Percentage of Patients Represented in HHAs with At Least 20 Valid Episodes, By Population Group

Population Group	# of Patients	% of Patients	
Dual Enrollment in Medicare	No	2,394,590	78.45%
and Medicaid	Yes	657,741	21.55%
Currently or Originally Eligible	No	2,405,389	78.80%
for Medicare due to Disability	Yes	646,942	21.20%
Location of HHA by Census	Northeast	631,656	20.69%
Region	Midwest	650,507	21.31%
	South	1,222,508	40.05%
	West	535,870	17.56%
	Missing	11,790	0.39%

# **1.7.** If there are differences in the data or sample used for different aspects of testing (e.g., reliability, validity, exclusions, risk adjustment), identify how the data or sample are different for each aspect of testing reported below.

To calculate the intra-class correlation (ICC) as part of reliability testing, the measure developer included Medicare-certified agencies with at least 40 home health quality episodes ending January 1, 2016 and December 31, 2016 and meeting the measure denominator criteria.<sup>1</sup> There were 7,862 such agencies (70.4 percent of the 11,166 agencies with at least one quality episode meeting the measure denominator criteria ending during the same time period). The sample included all quality episodes at these agencies (3,825,093 in total) meeting the measure denominator criteria ending between January 1, 2016 to December 31, 2016.

To calculate the beta-binomial scores (as part of reliability testing) and conduct analyses related to validity testing and exclusions, the measure developer included Medicare-certified agencies with at least 20 home health quality episodes ending between January 1, 2016 to December 31, 2016 and meeting the measure denominator criteria. There were 9,028 such agencies (80.9 percent of the 11,166 agencies with at least one quality episode meeting the measure denominator criteria ending during the same time period). The sample included all quality episodes meeting the measure denominator criteria at these agencies (3,858,808 in total) ending between January 1, 2016 to December 31, 2016.

The risk adjustment model was developed using OASIS national repository data from assessments submitted between January 1, 2016 and December 31, 2016 (which included ~ 6.4 million episodes of care).

1.8 What were the social risk factors that were available and analyzed? For example, patient-reported data (e.g., income, education, language), proxy variables when social risk data are not collected from each patient (e.g. census tract), or patient community characteristics (e.g. percent vacant housing, crime rate) which do not have to be a proxy for patient-level data.

We examined social risk factors that are available from the OASIS, as shown below. For operational and logistical reasons related to the monthly processing of this measure, drawing risk factors from external sources is not currently possible.

- Sex (female, male)
- Age (in 10 categories)
- Payment source (proxy for Medicaid coverage and dual eligibility using M0150 Current Payment Sources for Home Care see table below for the OASIS item responses).

Response for M0150 – Current Payment Sources for Home Care (Mark all that apply)

M0150	Responses
0	None; no charge for current service
1	Medicare (traditional fee-for-service)
2	Medicare (HMO/managed care/Advantage plan)

M0150	Responses
3	Medicaid (traditional fee-for-service)
4	Medicaid (HMO/managed care)
5	Workers' compensation
6	Title programs (for example, Title III, V, or XX)
7	Other government (for example, TriCare, VA)
8	Private insurance
9	Private HMO/managed care
10	Self-pay
11	Other (specify)
UK	Unknown

#### 2a2. RELIABILITY TESTING

<u>Note</u>: If accuracy/correctness (validity) of data elements was empirically tested, separate reliability testing of data elements is not required – in 2a2.1 check critical data elements; in 2a2.2 enter "see section 2b2 for validity testing of data elements"; and skip 2a2.3 and 2a2.4.

2a2.1. What level of reliability testing was conducted? (may be one or both levels)

Critical data elements used in the measure (*e.g., inter-abstractor reliability; data element reliability must address ALL critical data elements*)

☑ **Performance measure score** (e.g., *signal-to-noise analysis*)

# 2a2.2. For each level checked above, describe the method of reliability testing and what it tests (*describe the steps—do not just name a method; what type of error does it test; what statistical analysis was used*)

Below, we address reliability at two levels: (1) the performance measure and (2) the underlying data element - OASIS item M1242 (Frequency of pain interfering with patient's activity or movement).

• Reliability of the Performance Measure Score: Abt measured the extent to which differences in each quality measure were due to actual differences in agency performance versus variation that arises from measurement error. Statistically, reliability depends on performance variation for a measure across agencies, the random variation in performance for a measure within an agency's panel of attributed beneficiaries, and the number of beneficiaries attributed to the agency. High reliability for a measure suggests that comparisons of relative performance across agencies are likely to be stable over different performance periods, and that the performance of one agency on the quality measure can confidently be distinguished from another. Potential reliability values range from zero to one, where one (highest possible reliability) means that all variation in the measure's rates is the result of variation in differences in performance across agencies, while zero (lowest possible reliability) means that all variation is a result of measurement error.

Following the approach described by Adams,<sup>26</sup> Abt fit a beta-binomial model to estimate measure reliability. The beta-binomial model is appropriate because a particular agency's measure rate follows a binomial distribution (i.e., all measures are pass/fail), and it is reasonable to assume that the agencies' true measure rates vary and follow a beta distribution. It is reasonable to use the beta distribution to fit the true measure rates because it is a flexible distribution on the interval from 0 to 1, can have any mean on the interval, and can be skewed left, right, or U-shaped.

<sup>&</sup>lt;sup>26</sup> For more information about reliability testing for performance measurement, as well as the methodology for constructing the reliability score reported on Table 6, see "Reliability of Provider Profiling: A Tutorial" by John Adams, RAND. <u>http://www.rand.org/pubs/technical\_reports/TR653.html</u>

Equation (1), which is based on the beta-binomial model, shows that reliability is dependent on two variance components: the variation across agencies, and variation within agencies. In general, reliability for agencies will be higher when the measure rates across agencies are more heterogeneous (as measured by the agency-to-agency variation). Agencies with larger samples (n) and pass rates (p) nearer to 0 or 1 will have higher levels of reliability because the agency-specific error is reduced (i.e. the estimated agency rates are more precise).

 $Reliability = \frac{\sigma_{agency-to-agency}^2}{\sigma_{agency-to-agency}^2 + \sigma_{agency-specific-error}^2} = \frac{\sigma_{agency-to-agency}^2}{\sigma_{agency-to-agency}^2 + \frac{p(1-p)}{n}}$ (1)

Abt also calculated the test-retest reliability using the ICC to measure between-agency variation and within-agency variation. First, we randomly divided home health episodes within each agency into two separate equally-sized groups. Then, we calculated performance rates for each group. Then, using the paired performance rates, we calculated the statistics absolute-agreement ICC (AA-ICC or ICC(2,1)) and consistency-of-agreement ICC (CA-ICC or ICC(3,1)). ICC values that approach 1 indicate that the fraction of the total variance due to between-agency variation is high.

- <u>Reliability of the Underlying Data Element:</u> The measure is calculated by comparing patient functioning at the start and end of a home health quality episode, as reported by the home health OASIS-C2 data set. Pain interfering with activity is based on response to OASIS-C2 item M1242 (Frequency of pain interfering with patient's activity or movement):
  - 0 Patient has no pain.
  - 1 Patient has pain that does not interfere with activity or movement.
  - 2 Less often than daily.
  - 3 Daily, but not constantly.
  - 4 All of the time.

In 2016 and 2017, Abt and partners conducted a field test of new and existing OASIS items on 12 HHAs in four states for 213 home health patients.<sup>27</sup> Home health registered nurses and physical therapists, trained by the study team, collected data during home visits at start of care (SOC) or resumption of care (ROC), and/or at discharge. Follow-up visits were conducted within 24 hours of the initial field test visit, by a different registered nurse or physical therapist to test interrater reliability. M1242 was one of the existing OASIS-C2 items that was tested. Interrater reliability was assessed for SOC or ROC and at Discharge with a linear weighted kappa. The number patients for which inter-rater reliability could be tested was 105 at SOC/ROC and 84 at discharge.

The kappa statistic is generally considered to be the "gold standard" statistic associated with item reliability as it factors in the possibility of chance agreement. Kappa values are reported as decimal values between 0.00 (poor) and 1.00 (perfect). These can be interpreted using the following seven categories:<sup>28</sup>

- Poor < 0.10</p>
- Slight = 0.10 to 0.20
- Fair = 0.21 to 0.40
- Moderate = 0.41 to 0.60
- Substantial = 0.61 to 0.80
- Near perfect= 0.81 to 0.99
- Perfect = 1.00

<sup>&</sup>lt;sup>27</sup> Abt Associates (2018). "OASIS Field Test Summary Report: Outcome and Assessment Information Set (OASIS) Quality Measure Development and Maintenance Project."

<sup>&</sup>lt;sup>28</sup> Landis JR, Koch GG. The measurement of observer agreement for categorical data. Biometrics, 1977. 33(1):159-174.

## 2a2.3. For each level of testing checked above, what were the statistical results from reliability testing?

(e.g., percent agreement and kappa for the critical data elements; distribution of reliability statistics from a signal-to-noise analysis)

• **<u>Reliability of the Performance Measure Score</u>**: The table below summarizes the distribution of reliability scores for the 9,028 agencies that had at least 20 valid episodes.

Distribution of Beta Binomial Reliability Scores for Agencies with at Least 20 Valid Episodes

Mean	Minimum	10 <sup>th</sup> Percentile	25th Percentile	Median	75th Percentile	90th Percentile	Maximum
0.95	0.74	0.87	0.93	0.97	0.99	1.00	1.00

For agencies with at least 40 valid episodes (recall that an ICC statistic is derived from paired performance rates), the AA-ICC is **0.900**, and the CA-ICC is also **0.900**.

• <u>Reliability of the Underlying Data Element</u>: The inter-rater reliability (weighted kappa) values for M1242 (frequency of pain interfering with patient's activity or movement) was 0.53 at SOC/ROC and 0.45 at discharge.

2a2.4 What is your interpretation of the results in terms of demonstrating reliability? (i.e., what do the results mean and what are the norms for the test conducted?)

<u>Reliability of the Performance Measure Score</u>: Using the beta-binomial model, Abt concluded that the measure reliability was high. The mean and median reliability scores of 0.95 and 0.97, respectively, are above the range considered acceptable (0.70 – 0.80) for drawing inferences about home health agencies.

The ICC statistics also suggest acceptable test-retest reliability.

• <u>Reliability of the Underlying Data Element</u> The weighted kappa statistic for inter-rater reliability indicated moderate agreement at SOC/ROC (0.53) and EOC (0.45). We conclude that the item achieves sufficient reliability.

#### **2b1. VALIDITY TESTING**

2b1.1. What level of validity testing was conducted? (may be one or both levels)

Critical data elements (data element validity must address ALL critical data elements)

#### ⊠ Performance measure score

⊠ Empirical validity testing

□ Systematic assessment of face validity of <u>performance measure score</u> as an indicator of quality or resource use (*i.e., is an accurate reflection of performance on quality or resource use and can distinguish good from poor performance*) NOTE: Empirical validity testing is expected at time of maintenance review; if not possible, justification is required.

2b1.2. For each level of testing checked above, describe the method of validity testing and what it tests (describe the steps—do not just name a method; what was tested, e.g., accuracy of data elements compared to authoritative source, relationship to another measure as expected; what statistical analysis was used)

Below, we address validity at two levels: (1) the performance measure and (2) the underlying data element - OASIS item M1242 (frequency of pain interfering with patient's activity or movement).

• <u>Validity of the Performance Measure Score:</u> Abt assessed the convergent validity of the measure. Convergent validity refers to the extent to which measures that are designed to assess the same construct are related to each other. To evaluate the convergent validity of the measure, Abt calculated the Spearman rank correlations of the *Improvement in Pain Interfering with Activity*  measure with other relevant measures, including the publicly-reported measures of home health quality derived from OASIS assessments.

Abt also calculates and reports the Spearman rank correlation of the *Improvement in Pain Interfering with Activity* measure with the Quality of Patient Care Star Rating. The Spearman rank correlation assesses the statistical dependence between the rankings of two variables. In our case, we rank HHAs according to the *Improvement in Pain Interfering with Activity* measure and other OASIS-based measures. High correlation or association between the *Improvement in Pain Interfering with Activity* measure and other functional measures of improvement would be expected and desired. Low correlation would indicate that the measure may not be valid (is not measuring what we think it is measuring).

• <u>Validity of the Underlying Data Element</u>: The Pain Interfering with Activity item has been used continuously as part of the OASIS since 2001. The behaviorally benchmarked responses were updated and improved based on input from clinicians and technical experts. The OASIS instrument has been published in the Federal Register for comment (both items and measures based off those items) and no objections or suggestions for revision have been noted regarding the response options.

The original OASIS item was originally carefully designed for measuring and ultimately enhancing patient outcomes as part of the National OBQI Demonstration project (1995 – 2000). OASIS items were derived by first specifying a set of patient outcomes considered critical by home care experts (e.g., nurses, physicians, therapists, social workers, administrators) for evaluating the effectiveness of care. These outcomes were chosen from the most important domains of health status addressed by home care providers. OASIS data items were developed, tested in hundreds of agencies, and refined for measuring outcomes in order to evaluate and enhance the effectiveness of home care. OASIS data items and measurement methods were reviewed by multidisciplinary panels of research methodologists, clinicians, home care managers, and policy analysts. Several tests of validity were conducted for each OASIS item, including Pain Interfering with Activity. Validity testing included:

1) Consensus validity by expert researcher/clinical panels for outcome measurement and risk factor measurement

2) Consensus validity by expert clinical panels for patient assessment and care planning

3) Criterion or convergent/predictive validity for outcome measurement/risk factor measurement

4) Convergent/predictive validity: case mix adjustment for payment

5) Validation by patient assessment and care planning

Descriptions for these validation assessments are taken from the "Volume 4 : OASIS Chronicle and Recommendation" OASIS and Outcome-Based Quality Improvement in Home Health Care, November 2001, Center for Health Services Research, University of Colorado Health Sciences Center, Denver, CO.

#### **2b1.3.** What were the statistical results from validity testing? (e.g., correlation; t-test)

• <u>Validity of the Performance Measure Score</u>: The table below shows the Spearman rank correlations of the *Improvement in Pain Interfering with Activity* measure with other publicly-reported measures of home health quality derived from OASIS assessments.

# Spearman Rank Correlations of *Improvement in Pain Interfering with Activity* Measure with Other Measures of Home Health Quality

Home Health Quality Measures	Spearman Rank Correlations
Improvement in Ambulation/Locomotion	0.6163
Improvement in Bathing	0.6861
Improvement in Bed Transfer	0.5221
Improvement in Management of Oral Medications	0.5054
Quality of Patient Care Star Ratings	0.6546

#### • Validity of the Underlying Data Element: As noted above in 2b1.2,

- 1. *Consensus validity:* The item was reviewed by panels of researchers and clinicians and was recommended for measuring patient outcomes relevant to home health care provision and quality measurement, or for risk adjustment of outcome analyses.
- 2. Consensus validity by expert clinical panels for patient assessment and care planning: The item was reviewed by a panel of clinical experts and was recommended for inclusion in a core set of data items for patient assessment and care planning.
- 3. Criterion or convergent/predictive validity for outcome measurement/risk factor measurement: The item was tested empirically for use in conjunction with outcome measures or risk factors predictive of patient outcomes. The item was found to be related to other indicators of health status and patient outcomes in a statistically significant and clinically meaningful way.
- 4. *Convergent/predictive validity:* Case mix adjustment for payment: The item was tested and is used in the grouping algorithm that, in part, determines the per-episode payment to home health agencies for care provided under the Medicare home health benefit.
- 5. *Validation by patient assessment and care planning:* The item has been used by clinicians for patient assessment and care planning in several hundred home health agencies and has been reported by practicing clinicians to be effective and useful for these purposes.

Results of these validation assessments are taken from the "Volume 4 : OASIS Chronicle and Recommendation" OASIS and Outcome-Based Quality Improvement in Home Health Care, November 2001, Center for Health Services Research, University of Colorado Health Sciences Center, Denver, CO.

**2b1.4. What is your interpretation of the results in terms of demonstrating validity**? (i.e., what do the results mean and what are the norms for the test conducted?)

- <u>Validity of the Performance Measure Score:</u> As detailed in the Spearman Rank Correlations table, the *Improvement in Pain Interfering with Activity* measure displays a statistically significant positive correlation with several publicly-reported measures that similarly assess patient functioning and the quality of home health care, which lends evidence to the measure's validity. It may be that strong performance on the other measures directly leads to an improvement in pain interfering with activity. It may also be the case that high quality agencies perform well on both the *Improvement in Pain Interfering with Activity* measure and other OASIS-based measures of patient functioning and communication due to cultural or organization-level factors.
- <u>Validity of the Underlying Data Element:</u> Item validity was established based on results of testing described in section 2b2.2, above. In addition, the item was also reviewed as part of the OMB/PRA review process for the most recent OASIS data set revision which allowed for two national comment periods (60 days and 30 days) wherein the face validity of the item was supported by the comments received.

#### **2b2. EXCLUSIONS ANALYSIS**

#### NA 🗆 no exclusions— skip to section <u>2b3</u>

# 2b2.1. Describe the method of testing exclusions and what it tests (*describe the steps—do not just name a method; what was tested, e.g., whether exclusions affect overall performance scores; what statistical analysis was used*)

There are two major exclusion types for the *Improvement in Pain Interfering with Activity* measure, including exclusions that are applicable to home health measures in general (i.e., generic exclusions) and exclusions that are specific to the *Improvement in Pain Interfering with Activity* measure. Generic exclusions include (i) children and maternity patients and (ii) non-Medicare/non-Medicaid patients.

Exclusions that are specific to the *Improvement in Pain Interfering with Activity* measure include (i) episodes of care that did not end in discharge to community, (ii) episodes did not have pain interfering with activity at baseline, and (iii) episodes in which the patient was non-responsive at baseline and therefore not expected to improve in pain interfering with activity.

Abt calculated the frequency of the exclusions that are specific to the *Improvement in Pain Interfering with Activity*, by exclusion type.

# **2b2.2.** What were the statistical results from testing exclusions? (*include overall number and percentage of individuals excluded, frequency distribution of exclusions across measured entities, and impact on performance measure scores*)

Home Health Stays	# of Episodes Excluded	% of Episodes Excluded	# of Episodes Remaining
A. All home health episodes	N/A	N/A	6,437,455
B. Home health episodes that exclude episodes that did not end in discharge to community	1,764,228	27.4	4,673,227
C. Home health episodes from B that exclude episodes for which the patient, at start/resumption of care, did not have pain interfering with activity	773,309	16.5	3,899,918
D. Home health episodes from C that exclude episodes in which the patient is nonresponsive	23,566	0.6	3,876,352

#### Measure Denominator Exclusion, January 2016 to December 2016

**2b2.3.** What is your interpretation of the results in terms of demonstrating that exclusions are needed to prevent unfair distortion of performance results? (*i.e.*, the value outweighs the burden of increased data collection and analysis. <u>Note</u>: *If patient preference is an exclusion*, the measure must be specified so that the effect on the performance score is transparent, e.g., scores with and without exclusion)

All the measure exclusions are conceptually justified, so the measure developer did not conduct further statistical analyses to test the exclusions. The remainder of this response provides justifications for the exclusions for the *Improvement in Pain Interfering with Activity* measure.

Exclusions that are specific to the *Improvement in Pain Interfering with Activity* measure include (i) episodes of care that did not end in discharge to community (i.e., episodes of care that ended in transfer to inpatient facility or death at home), (ii) episodes in which the patient did not have pain interfering with activity at baseline, and (iii) episodes in which the patient was non-responsive at baseline. For exclusion (i), the information needed to calculate the measure is not collected for these episodes of care. Exclusions (ii), and

(iii) are justified because it would be impossible for these patients to demonstrate measurable improvement in pain interfering with activity over the episode of care.

The generic exclusions for this measure include:

- *Children And Maternity Patients* The OASIS data set items are designed to be collected for nonmaternity, adult patients who are 18 years and older. Maternity patients, and patients less than 18 years of age are excluded.
- Non-Medicare/non-Medicaid Patients Medicare-certified home health agencies are currently required to collect and submit OASIS data only on Medicare and Medicaid patients who are receiving skilled home health care. .

If the agency sample includes fewer than 20 episodes after all other patient-level exclusions are applied, or if the agency has been in operation less than six months, then the data is suppressed from public reporting on Home Health Compare.

#### 2b3. RISK ADJUSTMENT/STRATIFICATION FOR OUTCOME OR RESOURCE USE MEASURES

If not an intermediate or health outcome, or PRO-PM, or resource use measure, skip to section <u>2b4</u>.

2b3.1. What method of controlling for differences in case mix is used?

□ No risk adjustment or stratification

Statistical risk model with <u>114</u> risk factors

□ Stratification by \_risk categories

 $\Box$  Other,

**2b3.1.1** If using a statistical risk model, provide detailed risk model specifications, including the risk model method, risk factors, coefficients, equations, codes with descriptors, and definitions.

The Improvement in Pain Interfering with Activity risk adjustment model includes 114 risk factors. The specification of the risk factors, estimated coefficients, and methodology are provided in the attachment.

2b3.2. If an outcome or resource use component measure is <u>not risk adjusted or stratified</u>, provide <u>rationale</u> <u>and analyses</u> to demonstrate that controlling for differences in patient characteristics (case mix) is not needed to achieve fair comparisons across measured entities.

Not applicable; this measure is risk-adjusted.

2b3.3a. Describe the conceptual/clinical <u>and</u> statistical methods and criteria used to select patient factors (clinical factors or social risk factors) used in the statistical risk model or for stratification by risk (*e.g., potential factors identified in the literature and/or expert panel; regression analysis; statistical significance of p<0.10; correlation of x or higher; patient factors should be present at the start of care) Also discuss any "ordering" of risk factor inclusion; for example, are social risk factors added after all clinical factors?* 

The risk adjustment model was developed using OASIS national repository data from assessments submitted between January 1, 2016 and December 31, 2016 (~6.4 million episodes of care). The risk factors used in the unique prediction model created for each outcome measure are derived from OASIS data collected during the start of care or resumption of care assessment. The risk factors were developed and reviewed by home health clinicians. No ordering was used to determine risk factor inclusion, though, as described below, statistical criteria were applied to remove risk factors that were not statistically significant.

The risk adjustment methodology used is based on logistic regression analysis which results in a statistical prediction model for each outcome measure. For each home health agency patient who is included in the denominator of the outcome measure, the model is used to calculate the predicted probability that the patient will experience the outcome. The predicted probability for a patient is calculated using the following formula:

$$P(x) = 1/\left(1 + e^{-(a + \sum b_i x_i)}\right)$$

Where:

*P*(*x*) = predicted probability of achieving outcome x

a = constant parameter listed in the model documentation

 $b_i$  = coefficient for risk factor i in the model documentation

 $x_i$  = value of risk factor i for this patient

Predicted probabilities for all patients included in the measure denominator are then averaged to derive an expected outcome value for the agency. This expected value is then used, together with the observed (unadjusted) outcome value and the expected value for the national population of home health agency patients for the same data collection period, to calculate a risk-adjusted outcome value for the home health agency. The formula for the adjusted value of the outcome measure is as follows:

$$X(A_{ra}) = X(A_{obs}) + X(N_{exp}) - X(A_{exp})$$

Where:

X(A<sub>ra</sub>) = Agency risk-adjusted outcome measure value

X(A<sub>obs</sub>) = Agency observed outcome measure value

 $X(A_{exp})$  = Agency expected outcome measure value

 $X(N_{exp})$  = National expected outcome measure value

If the result of this calculation is a value greater than 100%, the adjusted value is set to 100%. Similarly, if the result is a negative number the adjusted value is set to zero.

**2b3.3b.** How was the conceptual model of how social risk impacts this outcome developed? Please check all that apply:

- Published literature
- $oxed{internal}$  Internal data analysis
- □ Other (please describe)

#### 2b3.4a. What were the statistical results of the analyses used to select risk factors?

We reviewed recent studies on accounting for sociodemographic status (SDS) conducted by the National Academies of Medicine (NAM), the Office of the Assistant Secretary for Planning and Evaluation (ASPE), and NQF.<sup>29</sup> These studies tested SDS factors such as dual eligibility, rurality, race/ethnicity, and disability. While most of these variables are available via CMS data sources, we were not currently able to use other data sources to risk adjust this measure due to the operational requirements of producing this measure on a monthly basis. However, in the future, we plan to further investigate using the CMS Enrollment Database and other geographic-level files (such as the Area Health Resource File or Census data) to incorporate these other factors into the risk adjustment model.

We therefore were only able to include variables available on the OASIS. These include gender, payment source, age and race/ethnicity. We did not include race/ethnicity since it was not recommended as a proxy for social risk from the previous studies noted above. The payment source risk factor serves as a proxy for dual

<sup>&</sup>lt;sup>29</sup> National Academies of Sciences, Engineering, and Medicine. (2016). Accounting for Social Risk Factors in Medicare Payment: Identifying Social Risk Factors. National Academies Press; Office of the Assistant Secretary for Planning and Evaluation (2016). Report to Congress: Social Risk Factors and Performance Under Medicare's Value-Based Purchasing Programs. United States Department of Health and Human Services; National Quality Forum (2016). Early Results of SES Trial Reveal Need for Better Data and SES Variables. Available at: <u>http://www.qualityforum.org/SES\_Trial\_Period\_Update.aspx</u>

eligibility and Medicaid coverage. It tends to underreport dual eligibility and Medicaid coverage, but, as shown below, is important in explaining measure performance.

2b3.4b. Describe the analyses and interpretation resulting in the decision to select social risk factors (*e.g.* prevalence of the factor across measured entities, empirical association with the outcome, contribution of unique variation in the outcome, assessment of between-unit effects and within-unit effects.) Also describe the impact of adjusting for social risk (or not) on providers at high or low extremes of risk.

We first present the observed measure performance between 2012 and 2016 stratified by each of the social risk factors. We note the greater increase in measure performance occurring between 2015 and 2016 (and to some extent between 2014 and 2015) may be related to the inception of the Quality of Patient Care Star Ratings and Home Health Value Based Purchasing. Both programs rely upon home health quality measures. The Quality of Patient Care Star Rating is a composite of a subset of measures reported on Home Health Compare, including *Improvement in Pain Interfering with Activity*. The Home Health Value Based Purchasing program uses home health quality measures to generate a score that is compared across HHAs within a state (for nine states) and, depending on relative performance, can negatively or positively affect home health claims payment.

Differences in episode-level observed measure performance by gender were small, though, on average, males performed better on the measure than females in every year from 2012 to 2016.

	2012	2013	2014	2015	2016
Male	68.3%	68.8%	68.9%	71.1%	75.3%
Female	66.5%	66.9%	67.1%	69.5%	74.0%

#### Average Episode-Level Observed Measure Performance over Time, by Gender

Average episode-level observed measure values also differed by age group. Younger patients performed worse. The relationships were steady over time, though greater improvement at all ages occurred between 2015 and 2016, likely for reasons stated above.

	2012	2013	2014	2015	2016
0-54	62.0%	62.1%	62.6%	64.8%	69.2%
55-60	61.8%	62.1%	62.5%	64.8%	69.1%
60-65	62.7%	63.3%	63.3%	65.8%	70.2%
65-70	67.7%	68.2%	68.3%	70.6%	75.1%
70-75	68.2%	68.8%	68.9%	71.1%	75.5%
75-80	68.5%	68.8%	69.2%	71.6%	75.8%
80-85	68.5%	69.0%	69.0%	71.5%	75.8%
85-90	68.4%	68.8%	69.1%	71.5%	75.7%
90-95	68.0%	68.5%	68.6%	71.2%	75.3%
95+	67.4%	68.0%	68.1%	70.6%	74.4%

#### Average Episode-Level Observed Measure Performance over Time, by Age Category

Average episode-level observed measure values were lowest for patients using Medicaid as a payment source. Patients who indicated Medicare only performed the best on the measure.

#### Average Episode-Level Observed Measure Performance over Time, by Payment Source

	2012	2013	2014	2015	2016
Medicare and					
Medicaid	66.2%	66.6%	66.6%	68.2%	71.1%
Medicaid only	59.6%	59.2%	61.2%	63.6%	67.6%
Medicare only	67.9%	68.5%	68.7%	70.6%	75.0%

The following table displays the relevant estimated coefficients from the logistic regression model of *Improvement in Pain Interfering with Activity* on a full set of OASIS-based risk factors (see Section 2.3.1.1). This table shows that male patients, older patients, and patients for whom the payer source is Medicare FFS are more likely to perform better on this measure. Almost all risk factors are statistically significant at the 1 percent statistical level.

	Coefficient	p-value
Female (excluded category)		
Male	0.104	0.000
AGE_0_54	-0.187	0.000
AGE_55_59	-0.215	0.000
AGE_60_64	-0.153	0.000
AGE_65_69 (excluded category)		
AGE_70_74	0.057	0.000
AGE_75_79	0.115	0.000
AGE_80_84	0.157	0.000
AGE_85_89	0.181	0.000
AGE_90_94	0.173	0.000
AGE_95PLUS	0.123	0.000
PAY_MCAID_ONLY	-0.192	0.000
PAY_MCARE _FFS (excluded		
category		
PAY_MCAREANDMCAID	-0.123	0.000
PAY_MCARE_HMO	-0.113	0.000
PAY_OTHER	-0.140	0.000

To address the second part of this question – regarding the impacts of not adjusting for certain social risk factors for providers at extreme levels of risk, we take a closer look at the HHA's geographic locations – specifically, we compare observed to risk adjusted measure values for HHAs located in rural versus urban settings. Rural residents may have worse health outcomes and experience reduced access to health services, affecting their ability to improve on this measure.<sup>30</sup> The table below shows observed and risk adjusted

<sup>&</sup>lt;sup>30</sup> Befort, C. A., Nazir, N., & Perri, M. G. (2012). Prevalence of obesity among adults from rural and urban areas of the United States: findings from NHANES (2005-2008). *The Journal of Rural Health*, *28*(4), 392-397.

Dye, C., Willoughby, D., Aybar-Damali, B., Grady, C., Oran, R., & Knudson, A. (2018). Improving Chronic Disease Self-Management by Older Home Health Patients through Community Health Coaching. *International journal of environmental research and public health*, *15*(4), 660.

measure values over time for rural and urban HHAs using a CBSA-based designation provided in the Provider of Services file

		2011	2012	2013	2014	2015	2016
Observe	Rural	59.5%	59.9%	59.3%	59.8%	61.5%	66.3%
d	Urban	64.7%	65.2%	65.6%	65.1%	66.4%	70.0%
Risk	Rural	61.4%	61.8%	61.2%	61.5%	63.0%	67.8%
Adjusted	Urban	65.3%	65.8%	66.4%	66.1%	67.6%	71.3%

Risk-Adjusted and Observed Improvement in Pain Interfering with Activity Measure Values by Rural and Urban Designation

This table shows that there are differences between rural and urban HHAs with this measure. Urban HHAs tended to perform better than rural HHAs. This was true for both observed and risk adjusted measure performance. These differences should be monitored going forward and consideration should be given to incorporating an external data source on rurality. As mentioned above, the greater increase in measure performance occurring between 2015 and 2016 (and to some extent between 2014 and 2015) may be related to the inception of the Quality of Patient Care Star Ratings and Home Health Value Based Purchasing.

**2b3.5.** Describe the method of testing/analysis used to develop and validate the adequacy of the statistical model <u>or</u> stratification approach (describe the steps—do not just name a method; what statistical analysis was used)

Using the assessment data from January 1, 2016 to December 31, 2016, nearly 6.4 million episodes of care were created. This was done by linking the start of care (SOC) or resumption of care (ROC) assessment for a patient with that patient's last assessment (i.e., transfer, discharge, or death). We split the population of 6.4 million episodes for calendar year 2016 in half such that 3.2 million episodes were used as a developmental sample and 3.2 million episodes were used as a validation sample. A structured approach was used to develop the initial prediction model. The risk factors used in the prediction models are derived from OASIS data collected during the start of care or resumption of care assessment. Because there were a large number of possible risk factors that needed to be considered for the measure outcome and because some of the risk factors used previously are expected to be removed as part of the transition to OASIS-D in January 2019, the following process was used to identify unique contributing risk factors to the prediction model:

1. We identified risk factors based on OASIS items that will remain following the OASIS-D transition. We examine the statistical properties of the items to specify risk factors (e.g., we grouped item response when there was low prevalence of certain responses). Team clinicians then reviewed all risk factors for clinical relevance and we re-defined or updated risk factors as necessary. We then divided these risk factors into 35 content focus groups (e.g., ICD9-based conditions). Where possible, we defined risk factors such that they flagged mutually exclusive subgroups within each content focus group. When modelling these risk factors, we use the risk factor flag indicating independence as our exclusion category.

2. We use a logistic regression specification to estimate coefficients among the full set of candidate risk factors. Those risk factors that are statistically significant at probability <0.001 are kept for further review.

5. The list of risk factors that achieved the probability<0.01 level were reviewed. If one response option level of an OASIS-D item was on the list, then risk factors representing the other response option levels of that OASIS-D item were added to the list. For example, if response option levels 1 and 2 for M1800 Grooming were statistically significant at probability<0.01 for a particular outcome, then response option level 3 for M1800 Grooming was added to the list.

6. A fixed logistic regression was computed on the list of risk factors that had achieved probability<0.001 and the risk factors that were added to the list because they were other response options for OASIS-D items represented on the list.

7. Goodness of fit statistics (R<sup>2</sup> and c-statistic) as well as bivariate correlations between the risk factor and the outcome were computed for how well the predicted values generated by the prediction model were related to the actual outcomes.

8. The initial model was reviewed by a team of at least three experienced home health clinicians. Each risk factor was reviewed for its clinical plausibility in being related to the outcome measure in the direction indicated by the coefficient in the prediction equation and its bivariate relationship. Risk factors that were not clinically plausible were identified for elimination.

9. The risk factors that were deemed not clinically plausible were removed from the prediction model and steps 6 and 7 in this process were repeated. The resulting logistic regression equation was designated as the prediction model for the outcome.

10. The prediction model was applied to the validation sample and goodness of fit statistics were computed. If these statistics were similar to the goodness of fit statistics computed with the development sample, the model become a final model. If the statistics were not similar, then alternative approaches to model building were considered.

*Provide the statistical results from testing the approach to controlling for differences in patient characteristics (case mix) below.* 

## If stratified, skip to <u>2b3.9</u>

#### **2b3.6.** Statistical Risk Model Discrimination Statistics (e.g., c-statistic, R-squared):

The c-statistic is the area under the Receiver Operating Characteristic curve. Intuitively, it is defined as follows: Let Y=1 denote outcome attainment, Y=0 denote nonattainment, and  $\hat{p}$  denote the predicted probability that Y=1. Enumerate all possible pairs of sample patients for whom Y=1 for the first patient and Y=0 for the second patient. C is the proportion of such pairs where  $\hat{p}$  for the patient with Y=1 is larger than  $\hat{p}$  for the patient with Y=0. The overall model development sample c-statistic is **0.656**. The overall model validation sample cstatistic is **0.657**.

Because the risk adjustment model uses a logistic specification, we report McFadden's R<sup>2</sup> to summarize model fit. The traditional R<sup>2</sup> value for linear specifications is the squared correlation between predicted and observed values for all patients in the developmental or validation samples. McFadden's R<sup>2</sup> is conceptually similar and compares the likelihood the full model to an intercept-only model. The overall model development sample R<sup>2</sup> is **0.053**. The overall model validation sample R<sup>2</sup> is **0.051**.

## **2b3.7.** Statistical Risk Model Calibration Statistics (e.g., Hosmer-Lemeshow statistic):

With a validation sample of over 3 million episodes, the Hosmer-Lemeshow test will reject the null assumption of equality even if differences in average performance are small. As such, we prefer a visual inspection of the risk decile plot below, which compares the average predicted performance against the average observed performance for *Improvement in Pain Interfering with Activity*. The plot below shows that the predicted and observed values are similar and monotonically increasing with predicted probability, both of which indicate a well calibrated model. Additionally, we consider the R<sup>2</sup> statistics (included in response to **2b3.6**) to be sufficient indicators of model fit.

#### 2b3.8. Statistical Risk Model Calibration – Risk decile plots or calibration curves:



Risk Deciles, Improvement in Pain Frequency Calendar Year 2016 Episodes

#### 2b3.9. Results of Risk Stratification Analysis:

Not applicable.

**2b3.10.** What is your interpretation of the results in terms of demonstrating adequacy of controlling for differences in patient characteristics (case mix)? (i.e., what do the results mean and what are the norms for the test conducted)

The c-statistic for the development sample is **0.656**, which is similar to the validation sample value of **0.657**, showing that the model differentiates between outcomes as well on new data as it does on the development data.

The McFadden's R<sup>2</sup> for the development sample is **0.053**, which is similar to the validation sample value of **0.051**, showing that the model is capable of describing the relationship between the covariates and the outcome in the development data set while also successfully predicting the outcome on a new data set.

**2b3.11.** Optional Additional Testing for Risk Adjustment (*not required*, but would provide additional support of adequacy of risk model, e.g., testing of risk model in another data set; sensitivity analysis for missing data; other methods that were assessed)

None

#### 2b4. IDENTIFICATION OF STATISTICALLY SIGNIFICANT & MEANINGFUL DIFFERENCES IN PERFORMANCE

2b4.1. Describe the method for determining if statistically significant and clinically/practically meaningful differences in performance measure scores among the measured entities can be identified (*describe the steps*—do not just name a method; what statistical analysis was used? Do not just repeat the information provided related to performance gap in 1b)

To demonstrate that the *Improvement in Pain Interfering* with Activity measure exhibits variation and that the variation is meaningful in discriminating performance among home health agencies, we conducted the following analyses:

- 1. First, we show that there is variation in the measure by examining the measure distribution mean, median, 10<sup>th</sup>, 25th, 75<sup>th</sup>, and 90<sup>th</sup> percentile values. We also calculated the truncated coefficient of variation (TCV).
  - a. We show that the measure is not "topped-out;" that is, we show there is room for improvement in the measure. Measures that are "topped-out" or close to being so are less able to meaningfully discriminate between providers. That is, if the majority of agencies are already performing at a high level, the measure is less able to distinguish between providers. We demonstrate that the 10<sup>th</sup> percentile value of the measure is less than 70 percent. That is, if the HHAs performing at the 10<sup>th</sup> percentile had a measure value of 70 percent, then we would consider the measure having little room for improvement.
  - b. We show that the interquartile range (IQR) is substantial. The IQR is calculated by subtracting the 25<sup>th</sup> percentile measure value from the 75<sup>th</sup> percentile measure; it shows the measure "spread."
  - c. The TCV is another measure of variation it is the ratio of the truncated standard deviation and truncated mean. We truncate by removing the bottom 5<sup>th</sup> percentile and the top 95<sup>th</sup> percentile of HHAs. A larger TCV indicates higher variability of the measure.
  - d. We show the same information for HHAs stratified by whether the census region in which the HHA is located.
- 2. Demonstrating that there is variation in the measure is not sufficient for concluding that the variation is meaningful. To examine whether the measure is meaningful in distinguishing performance across agencies, we examined the performance of the measure by an altered version of the Quality of Patient Care (QoPC) Star Rating and tested whether measure values differ by rating and whether the difference is statistically significant at the 5 percent significance level. The QoPC Star Rating is composed of eight equally weighted quality measures, including *Improvement in Pain Interfering with Activity*.<sup>31</sup> We created an altered version that removes the *Improvement in Pain Interfering with Activity* from the QoPC Star Ratings (keeping the remaining measures and methodology the same). The other measures include other functional improvement measures, two process measures and a claims-based hospitalization measure. The QoPC Star Ratings are a composite of these measures and take on nine values (1 to 5 stars in half star increments). Higher stars indicate higher quality. We thus expect that HHAs with higher QoPC Star Ratings (or alternate) values will have higher values on the *Improvement in Pain Interfering with Activity* measure.

**2b4.2.** What were the statistical results from testing the ability to identify statistically significant and/or clinically/practically meaningful differences in performance measure scores across measured entities? (e.g., number and percentage of entities with scores that were statistically significantly different from mean or some benchmark, different from expected; how was meaningful difference defined)

The table below shows the distribution of the *Improvement in Pain Interfering with Activity* measure across the 9,028 agencies that had at least 20 episodes available. The median is 70.0 percent. The 10<sup>th</sup> percentile value is 31.6 percent and the 90<sup>th</sup> percentile value is 92.9 percent. The IQR is 27.6 percent. The TCV (not shown in the table) is 36.2 percent. These statistics show that the measure is not topped out and there is still sufficient room for improvement.

<sup>&</sup>lt;sup>31</sup> <u>https://www.cms.gov/Medicare/Quality-Initiatives-Patient-Assessment-</u> Instruments/HomeHealthQualityInits/Downloads/QoPC-Methodology\_for\_April\_2018.pdf

#### Distribution of Improvement in Pain Interfering with Activity (Risk Adjusted) Overall and by Census Region

	#HHAs	Mean	10 <sup>th</sup>	25 <sup>th</sup>	50 <sup>th</sup>	75 <sup>th</sup>	90 <sup>th</sup>	IQR
All*	9,028	65.6%	31.6%	54.5%	70.0%	82.1%	92.9%	27.6%
Northeast	767	66.1%	37.5%	57.7%	70.4%	78.6%	87.5%	20.9%
Midwest	2,397	64.6%	33.3%	52.8%	68.0%	81.1%	91.4%	28.3%
South	4,017	63.0%	23.3%	50.0%	68.6%	81.0%	92.0%	31.0%
West	1,807	73.1%	50.0%	63.0%	74.4%	86.7%	96.0%	23.6%

\*Note that "All" includes all HHAs in the 50 states and U.S. territories. The census regions only include U.S. States (thus, the number of HHAs in each census region does not all up to "All").

This figure and table below shows the measure value by "altered" QoPC Star Rating. The figure shows that the *Improvement in Pain Interfering with Activity* measure steadily increases with a higher rating. The table below the figure shows the same information in table format. It includes the count of the number of HHAs with each rating as well as the statistical significance of a t-test between with sequential pairing. For example a t-test of the difference between the measure value for HHAs with 1.0 stars versus HHAs with 1.5 stars showed that the difference was different from zero with a p-value of 0.006 (i.e., statistically significant at the 5 percent level). All sequential pairwise differences were statistically significantly different from zero.



Measure Performance by "Altered" Quality of Patient Care Star Rating\*

		<b>Risk Adjusted Measure</b>	
Altered QoPC Star Rating	HHA Count	Value	Pairwise p-value
1.0	28	33.1%	-
1.5	255	42.9%	0.006
2.0	841	52.2%	0.000
2.5	1,335	63.0%	0.000
3.0	1,733	69.5%	0.000
3.5	1,835	74.5%	0.000
4.0	1,543	80.0%	0.000
4.5	938	85.8%	0.000
5.0	349	90.3%	0.000
Missing	171	64.4%	-

\*The QoPC Star Rating was altered by removing the Improvement in Pain Interfering with Activity measure from the rating calculation.

# **2b4.3.** What is your interpretation of the results in terms of demonstrating the ability to identify statistically significant and/or clinically/practically meaningful differences in performance across measured entities? (i.e., what do the results mean in terms of statistical and meaningful differences?)

Based on these findings, we conclude that the *Improvement in Pain Interfering with Activity* measure is able to produce meaningful differences across HHAs. First, the measure exhibits sufficient variation – it is not topped out and there is room for measure improvement among the majority of HHAs. Second, measure performance is related to other metrics in the direction expected with statistically significant differences in measure performance across strata.

#### 2b5. COMPARABILITY OF PERFORMANCE SCORES WHEN MORE THAN ONE SET OF SPECIFICATIONS

If only one set of specifications, this section can be skipped.

<u>Note</u>: This item is directed to measures that are risk-adjusted (with or without social risk factors) OR to measures with more than one set of specifications/instructions (e.g., one set of specifications for how to identify and compute the measure from medical record abstraction and a different set of specifications for claims or eMeasures). It does not apply to measures that use more than one source of data in one set of specification for the numerator). Comparability is not required when comparing performance scores with and without social risk factors in the risk adjustment model. However, if comparability is not demonstrated for measures with more than one set of specifications/instructions, the different specifications (e.g., for medical records vs. claims) should be submitted as separate measures.

2b5.1. Describe the method of testing conducted to compare performance scores for the same entities across the different data sources/specifications (*describe the steps—do not just name a method; what statistical analysis was used*)

N/A - one set of data/specifications are used

2b5.2. What were the statistical results from testing comparability of performance scores for the same entities when using different data sources/specifications? (*e.g., correlation, rank order*)

N/A - one set of data/specifications are used

2b5.3. What is your interpretation of the results in terms of the differences in performance measure scores for the same entities across the different data sources/specifications? (i.e., what do the results mean and what are the norms for the test conducted)

N/A - one set of data/specifications are used

#### 2b6. MISSING DATA ANALYSIS AND MINIMIZING BIAS

2b6.1. Describe the method of testing conducted to identify the extent and distribution of missing data (or nonresponse) and demonstrate that performance results are not biased due to systematic missing data (or differences between responders and nonresponders) and how the specified handling of missing data minimizes bias (*describe the steps—do not just name a method; what statistical analysis was used*)

There are minimal issues with missing data because the OASIS submission system rejects assessments with missing values. The provider must then resubmit the assessment.

2b6.2. What is the overall frequency of missing data, the distribution of missing data across providers, and the results from testing related to missing data? (*e.g., results of sensitivity analysis of the effect of various rules for missing data/nonresponse; if no empirical sensitivity analysis, identify the approaches for handling missing data that were considered and pros and cons of each*)

There are minimal issues with missing data because the OASIS submission system rejects assessments with missing values. The provider must then resubmit the assessment.

2b6.3. What is your interpretation of the results in terms of demonstrating that performance results are not biased due to systematic missing data (or differences between responders and nonresponders) and how the specified handling of missing data minimizes bias? (i.e., what do the results mean in terms of supporting the selected approach for missing data and what are the norms for the test conducted; if no empirical analysis, provide rationale for the selected approach for missing data)

There are minimal issues with missing data because the OASIS submission system rejects assessments with missing values. The provider must then resubmit the assessment.

# 3. Feasibility

Extent to which the specifications including measure logic, require data that are readily available or could be captured without undue burden and can be implemented for performance measurement.

#### **3a. Byproduct of Care Processes**

For clinical measures, the required data elements are routinely generated and used during care delivery (e.g., blood pressure, lab test, diagnosis, medication order).

#### 3a.1. Data Elements Generated as Byproduct of Care Processes.

Generated or collected by and used by healthcare personnel during the provision of care (e.g., blood pressure, lab value, diagnosis, depression score)

#### If other:

#### **3b. Electronic Sources**

The required data elements are available in electronic health records or other electronic sources. If the required data are not in electronic health records or existing electronic sources, a credible, near-term path to electronic collection is specified.

**3b.1.** To what extent are the specified data elements available electronically in defined fields (*i.e.*, data elements that are needed to compute the performance measure score are in defined, computer-readable fields) Update this field for maintenance of endorsement.

ALL data elements are in defined fields in electronic clinical data (e.g., clinical registry, nursing home MDS, home health OASIS)

**3b.2.** If ALL the data elements needed to compute the performance measure score are not from electronic sources, specify a credible, near-term path to electronic capture, OR provide a rationale for using other than electronic sources. For <u>maintenance of endorsement</u>, if this measure is not an eMeasure (eCQM), please describe any efforts to develop an eMeasure (eCQM).

**3b.3.** If this is an eMeasure, provide a summary of the feasibility assessment in an attached file or make available at a measure-specific URL. Please also complete and attach the NQF Feasibility Score Card.

#### Attachment:

#### **3c. Data Collection Strategy**

Demonstration that the data collection strategy (e.g., source, timing, frequency, sampling, patient confidentiality, costs associated with fees/licensing of proprietary measures) can be implemented (e.g., already in operational use, or testing demonstrates that it is ready to put into operational use). For eMeasures, a feasibility assessment addresses the data elements and measure logic and demonstrates the eMeasure can be implemented or feasibility concerns can be adequately addressed.

3c.1. <u>Required for maintenance of endorsement</u>. Describe difficulties (as a result of testing and/or operational use of the measure) regarding data collection, availability of data, missing data, timing and

frequency of data collection, sampling, patient confidentiality, time and cost of data collection, other feasibility/implementation issues.

<u>IF instrument-based</u>, consider implications for both individuals providing data (patients, service recipients, respondents) and those whose performance is being measured.

OASIS data collection and transmission is a requirement for the Medicare Home Health Conditions of Participation. Information on pain interfering with activity used to calculate this measure is recorded in the relevant OASIS items embedded in the agency's clinical assessment as part of normal clinical practice. OASIS data are collected by the home health agency during the care episode and transmitted electronically to the CMS national OASIS repository. No issues regarding availability of data, missing data, timing or frequency of data collection, patient confidentiality, time or cost of data collection, feasibility or implementation have become apparent since OASIS-C was implemented 1/1/2010.

**3c.2.** Describe any fees, licensing, or other requirements to use any aspect of the measure as specified (*e.g.,* value/code set, risk model, programming code, algorithm).

Not Applicable

# 4. Usability and Use

Extent to which potential audiences (e.g., consumers, purchasers, providers, policy makers) are using or could use performance results for both accountability and performance improvement to achieve the goal of highquality, efficient healthcare for individuals or populations.

#### 4a. Accountability and Transparency

Performance results are used in at least one accountability application within three years after initial endorsement and are publicly reported within six years after initial endorsement (or the data on performance results are available). If not in use at the time of initial endorsement, then a credible plan for implementation within the specified timeframes is provided.

#### 4.1. Current and Planned Use

NQF-endorsed measures are expected to be used in at least one accountability application within 3 years and publicly reported within 6 years of initial endorsement in addition to performance improvement.

Specific Plan for Use	Current Use (for current use provide URL)
	Public Reporting
	Home Health Compare
	http://www.cms.gov/HomeHealthCompare/search.aspx
	Home Health Compare
	http://www.cms.gov/HomeHealthCompare/search.aspx

#### 4a1.1 For each CURRENT use, checked above (update for <u>maintenance of endorsement</u>), provide:

- Name of program and sponsor
- Purpose
- Geographic area and number and percentage of accountable entities and patients included
- Level of measurement and setting

The Home Health Compare website is federal government website managed by the Centers for Medicare & Medicaid Services (CMS). It provides information to consumers about the quality of care provided by Medicare-certified home health agencies throughout the nation. The measures reported on Home Health Compare includes all Medicare-certified agencies with at least 20 home health quality episodes.

In the 12-month period ending December 31, 2016, there were 9,028 such agencies (80.9 percent of the 11,166 agencies with at least one quality episode) that met the measure denominator criteria for reporting of Improvement in Pain Interfering with Activity. This included 3,858,808 episodes of care nationally. CMS's Home Health Quality Initiative "Outcome Quality Measure Report" provides all Medicare-certified home health agencies with opportunities to use outcome measures for outcome-based quality improvement. The report allows agencies to benchmark their performance against other agencies across the state and nationally, as well as their own performance from prior time periods. All Medicare-certified home health agencies can access their Outcome Quality Measure Reportsvia CMS's online CASPER system.

**4a1.2.** If not currently publicly reported OR used in at least one other accountability application (e.g., payment program, certification, licensing) what are the reasons? (e.g., Do policies or actions of the developer/steward or accountable entities restrict access to performance results or impede implementation?) Not applicable

4a1.3. If not currently publicly reported OR used in at least one other accountability application, provide a credible plan for implementation within the expected timeframes -- any accountability application within 3 years and publicly reported within 6 years of initial endorsement. (*Credible plan includes the specific program, purpose, intended audience, and timeline for implementing the measure within the specified timeframes. A plan for accountability applications addresses mechanisms for data aggregation and reporting.*)

#### Not Applicable

4a2.1.1. Describe how performance results, data, and assistance with interpretation have been provided to those being measured or other users during development or implementation.

How many and which types of measured entities and/or others were included? If only a sample of measured entities were included, describe the full population and how the sample was selected.

All home health agencies with at least 20 qualifying episodes receive quarterly measure reports on all of their publicly-reported measures. In addition, providers can run on-demand, confidential reports showing individual measure results and national averages, through CMS' CASPER system. There is an email box that HHAs may submit questions to as well as a website on which the latest measure updates are posted.

# 4a2.1.2. Describe the process(es) involved, including when/how often results were provided, what data were provided, what educational/explanatory efforts were made, etc.

All home health agencies with at least 20 qualifying episodes receive quarterly measure reports on all of their publicly-reported measures. In addition, providers can run on-demand, confidential reports showing individual measure results and national averages, through CMS' CASPER system. There is an email box that HHAs may submit questions to as well as a website on which the latest measure updates are posted. The OASIS Guidance Manual describes the OASIS-based reports that are available as well as the sources of information for the reports. Instructions on using the reports for quality monitoring are provided, illustrated with sample reports from a hypothetical home care agency. It is designed to help home health agencies make use of the reports for monitoring and improving quality of care. Additionally, home health quality reporting program training was held in 2017.

# 4a2.2.1. Summarize the feedback on measure performance and implementation from the measured entities and others described in 4d.1.

#### Describe how feedback was obtained.

Home health agencies receive quarterly measure reports on all of their measures. There is an email box that HHAs may submit questions to as well as a website on which the latest measure updates are posted. Because of the changes made to the OASIS in OASIS D (effectively January 1, 2019), risk models for publically reported outcome measures have been updated. CMS will make available information about risk models and covariates on the website and the updated models will be available soon.

4a2.2.2. Summarize the feedback obtained from those being measured.

There is an email box that HHAs may submit regarding quality measures; all questions and responses are captured in an Access database for analysis and CMS receives quarterly reports on questions submitted. Thematic issues arising from the mailbox inform guidance to providers. As in 4a2.2.1.

#### 4a2.2.3. Summarize the feedback obtained from other users

There haven't been any requests for measure modification, nor any modifications made.

4a2.3. Describe how the feedback described in 4a2.2.1 has been considered when developing or revising the measure specifications or implementation, including whether the measure was modified and why or why not.

#### Not applicable for this time period.

#### Improvement

Progress toward achieving the goal of high-quality, efficient healthcare for individuals or populations is demonstrated. If not in use for performance improvement at the time of initial endorsement, then a credible rationale describes how the performance results could be used to further the goal of high-quality, efficient healthcare for individuals or populations.

4b1. Refer to data provided in 1b but do not repeat here. Discuss any progress on improvement (trends in performance results, number and percentage of people receiving high-quality healthcare; Geographic area and number and percentage of accountable entities and patients included.)

If no improvement was demonstrated, what are the reasons? If not in use for performance improvement at the time of initial endorsement, provide a credible rationale that describes how the performance results could be used to further the goal of high-quality, efficient healthcare for individuals or populations.

Tables 4 and 5 in the "Importance to Report" attachment show observed and predicted measure performance for population groups, respectively. For all population groups, measure performance has improved over time. The greatest improvement in measure performance between 2013 and 2016 for each population subgroup was for:

- Females (66.9 percent in 2013 to 74.0 percent in 2016)
- Whites (67.1 percent in 2013 to 74.8 percent in 2016)
- Under 65 (62.5 percent in 2013 to 69.5 percent in 2016)
- Disabled (63.2 percent in 2013 to 70.3 percent in 2016 similar to not disabled)
- Not dual (68.0 percent in 2013 to 75.2 percent in 2016 similar to dual)
- Large HHAs (53.2 percent in 2013 to 63.4 percent in 2016)
- HHAs in the South (66.2 percent in 2013 to 74.5 percent in 2016)

The subgroup with the smallest improvement in performance during this time period was for patients served by small HHAs (bottom 25th percentile in size). Performance for this subgroup only improved from 55.0 percent in 2013 to 56.2 percent in 2016. Note that the number of episodes for small HHAs was only 31,332 during 2016 (or 0.81 percent of episodes for which this measure is available).

There was generally fairly large improvement in measure performance during the 2013 to 2016 period. Overall, improvement was 6.9 percentage points and most population subgroups saw this level of improvement. The largest improvement occurred from 2015 to 2016 – more than half (4.4 percentage points) of the 2013-2016 improvement occurred between 2015 and 2016. We expect to see a similar phenomenon between 2016 and later years. This is likely due to the introduction of several initiatives that incorporate this measure – the Quality of Patient Care (QoPC) Star Ratings, a composite of this measure and several others that has been publicly reported on Home Health Compare since July 2015 and Home Health Value Based Purchasing (HHVBP). HHVBP began in 2016 and involves nine states. Several participating states encompass a large number of HHAs and providers in other states may be anticipating the expansion of this model.

#### 4b2. Unintended Consequences

The benefits of the performance measure in facilitating progress toward achieving high-quality, efficient healthcare for individuals or populations outweigh evidence of unintended negative consequences to individuals or populations (if such evidence exists).

# 4b2.1. Please explain any unexpected findings (positive or negative) during implementation of this measure including unintended impacts on patients.

Recent improvement in this measure has been relatively large compared to historical trends. We believe these large improvements are due to the implementation of two initiatives that involve this measure – the QoPC Star Ratings and HHVBP – beginning in 2015 and 2016.

#### 4b2.2. Please explain any unexpected benefits from implementation of this measure.

We do not report any unexpected benefits from implementation of this measure at this time.

# 5. Comparison to Related or Competing Measures

If a measure meets the above criteria <u>and</u> there are endorsed or new related measures (either the same measure focus or the same target population) or competing measures (both the same measure focus and the same target population), the measures are compared to address harmonization and/or selection of the best measure.

#### 5. Relation to Other NQF-endorsed Measures

Are there related measures (conceptually, either same measure focus or target population) or competing measures (conceptually both the same measure focus and same target population)? If yes, list the NQF # and title of all related and/or competing measures.

No

#### 5.1a. List of related or competing measures (selected from NQF-endorsed measures)

#### 5.1b. If related or competing measures are not NQF endorsed please indicate measure title and steward.

5a. Harmonization of Related Measures

The measure specifications are harmonized with related measures;

OR

The differences in specifications are justified

**5a.1.** If this measure conceptually addresses EITHER the same measure focus OR the same target population as NQF-endorsed measure(s):

Are the measure specifications harmonized to the extent possible?

No

5a.2. If the measure specifications are not completely harmonized, identify the differences, rationale, and impact on interpretability and data collection burden.

see 5b.1.

#### **5b.** Competing Measures

The measure is superior to competing measures (e.g., is a more valid or efficient way to measure); **OR** 

Multiple measures are justified.

5b.1. If this measure conceptually addresses both the same measure focus and the same target population as NQF-endorsed measure(s):

Describe why this measure is superior to competing measures (e.g., a more valid or efficient way to measure quality); OR provide a rationale for the additive value of endorsing an additional measure. (Provide analyses when possible.)

A search using the NQF QPS for outcome measures reporting rates of improvement in pain identified two measures used in the hospice setting (NQF# 0676, 0677 - Percent of Residents Who Self-Report Moderate to

Severe Pain). These measures are focused on inpatient (not homebound) patients, are calculated using data that are not currently collected in the home health setting, and do not consider the functional impact of pain.

# **Appendix**

**A.1 Supplemental materials may be provided in an appendix.** All supplemental materials (such as data collection instrument or methodology reports) should be organized in one file with a table of contents or bookmarks. If material pertains to a specific submission form number, that should be indicated. Requested information should be provided in the submission form and required attachments. There is no guarantee that supplemental materials will be reviewed.

Attachment Attachment: 0177\_Pain\_Importance\_to\_Report\_Tables.docx

# **Contact Information**

Co.1 Measure Steward (Intellectual Property Owner): Centers for Medicare & Medicaid Services

Co.2 Point of Contact: Joan, Proctor, Joan.Proctor2@cms.hhs.gov, 443-526-6938-

Co.3 Measure Developer if different from Measure Steward: Centers for Medicare & Medicaid Services

Co.4 Point of Contact: Joan, Proctor, Joan.Proctor2@cms.hhs.gov, 443-526-6938-

# **Additional Information**

Ad.1 Workgroup/Expert Panel involved in measure development

Provide a list of sponsoring organizations and workgroup/panel members' names and organizations. Describe the members' role in measure development.

#### Not Applicable

Measure Developer/Steward Updates and Ongoing Maintenance Ad.2 Year the measure was first released: 2004 Ad.3 Month and Year of most recent revision: 11, 2009 Ad.4 What is your frequency for review/update of this measure? Annual Ad.5 When is the next scheduled review/update for this measure? 09, 2018 Ad.6 Copyright statement: NA Ad.7 Disclaimers: NA Ad.8 Additional Information/Comments: NA